# Introduction to Econometrics

## THIRD EDITION UPDATE

**James H. Stock**
Harvard University

**Mark W. Watson**
Princeton University

# 18

# The Theory of Multiple Regression

This chapter provides an introduction to the theory of multiple regression analysis. The chapter has four objectives. The first is to present the multiple regression model in matrix form, which leads to compact formulas for the OLS estimator and test statistics. The second objective is to characterize the sampling distribution of the OLS estimator, both in large samples (using asymptotic theory) and in small samples (if the errors are homoskedastic and normally distributed). The third objective is to study the theory of efficient estimation of the coefficients of the multiple regression model and to describe generalized least squares (GLS), a method for estimating the regression coefficients efficiently when the errors are heteroskedastic and/or correlated across observations. The fourth objective is to provide a concise treatment of the asymptotic distribution theory of instrumental variables (IV) regression in the linear model, including an introduction to generalized method of moments (GMM) estimation in the linear IV regression model with heteroskedastic errors.

The chapter begins by laying out the multiple regression model and the OLS estimator in matrix form in Section 18.1. This section also presents the extended least squares assumptions for the multiple regression model. The first four of these assumptions are the same as the least squares assumptions of Key Concept 6.4 and underlie the asymptotic distributions used to justify the procedures described in Chapters 6 and 7. The remaining two extended least squares assumptions are stronger and permit us to explore in more detail the theoretical properties of the OLS estimator in the multiple regression model.

The next three sections examine the sampling distribution of the OLS estimator and test statistics. Section 18.2 presents the asymptotic distributions of the OLS estimator and $t$-statistic under the least squares assumptions of Key Concept 6.4. Section 18.3 unifies and generalizes the tests of hypotheses involving multiple coefficients presented in Sections 7.2 and 7.3, and provides the asymptotic distribution of the resulting $F$-statistic. In Section 18.4, we examine the exact sampling distributions of the OLS estimator and test statistics in the special case that the errors are homoskedastic and normally distributed. Although the assumption of homoskedastic normal errors is implausible in most econometric applications, the exact sampling distributions are of theoretical interest, and $p$-values computed using these distributions often appear in the output of regression software.

The next two sections turn to the theory of efficient estimation of the coefficients of the multiple regression model. Section 18.5 generalizes the Gauss–Markov theorem to multiple regression. Section 18.6 develops the method of generalized least squares (GLS).

The final section takes up IV estimation in the general IV regression model when the instruments are valid and strong. This section derives the asymptotic distribution of the TSLS estimator when the errors are heteroskedastic and provides expressions for the standard error of the TSLS estimator. The TSLS estimator is one of many possible GMM estimators, and this section provides an introduction to GMM estimation in the linear IV regression model. It is shown that the TSLS estimator is the efficient GMM estimator if the errors are homoskedastic.

***Mathematical prerequisite.*** The treatment of the linear model in this chapter uses matrix notation and the basic tools of linear algebra and assumes that the reader has taken an introductory course in linear algebra. Appendix 18.1 reviews vectors, matrices, and the matrix operations used in this chapter. In addition, multivariate calculus is used in Section 18.1 to derive the OLS estimator.

## 18.1 The Linear Multiple Regression Model and OLS Estimator in Matrix Form

The linear multiple regression model and the OLS estimator can each be represented compactly using matrix notation.

### The Multiple Regression Model in Matrix Notation

The population multiple regression model (Key Concept 6.2) is

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_k X_{ki} + u_i, i = 1, \ldots, n. \qquad (18.1)$$

To write the multiple regression model in matrix form, define the following vectors and matrices:

$$\boldsymbol{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \boldsymbol{U} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}, \boldsymbol{X} = \begin{pmatrix} 1 & X_{11} & \cdots & X_{k1} \\ 1 & X_{12} & \cdots & X_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{1n} & \cdots & X_{kn} \end{pmatrix} = \begin{pmatrix} \boldsymbol{X}_1' \\ \boldsymbol{X}_2' \\ \vdots \\ \boldsymbol{X}_n' \end{pmatrix}, \text{and } \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}, \quad (18.2)$$

so $Y$ is $n \times 1$, $X$ is $n \times (k + 1)$, $U$ is $n \times 1$, and $\boldsymbol{\beta}$ is $(k + 1) \times 1$. Throughout we denote matrices and vectors by bold type. In this notation,

- $Y$ is the $n \times 1$ dimensional vector of $n$ observations on the dependent variable.

- $X$ is the $n \times (k + 1)$ dimensional matrix of $n$ observations on the $k + 1$ regressors (including the "constant" regressor for the intercept).

- The $(k + 1) \times 1$ dimensional column vector $X_i$ is the $i^{\text{th}}$ observation on the $k + 1$ regressors; that is, $X_i' = (1 \ X_{1i} \ldots X_{ki})$, where $X_i'$ denotes the transpose of $X_i$.

- $U$ is the $n \times 1$ dimensional vector of the $n$ error terms.

- $\boldsymbol{\beta}$ is the $(k + 1) \times 1$ dimensional vector of the $k + 1$ unknown regression coefficients.

The multiple regression model in Equation (18.1) for the $i^{\text{th}}$ observation, written using the vectors $\boldsymbol{\beta}$ and $X_i$, is

$$Y_i = X_i'\boldsymbol{\beta} + u_i, i = 1, \ldots, n. \tag{18.3}$$

## The Extended Least Squares Assumptions in the Multiple Regression Model

KEY CONCEPT

18.1

The linear regression model with multiple regressors is

$$Y_i = X_i'\boldsymbol{\beta} + u_i, i = 1, \ldots, n. \tag{18.4}$$

The extended least squares assumptions are

1. $E(u_i | X_i) = 0$ ($u_i$ has conditional mean zero);
2. $(X_i, Y_i), i = 1, \ldots, n$, are independently and identically distributed (i.i.d.) draws from their joint distribution;
3. $X_i$ and $u_i$ have nonzero finite fourth moments;
4. $X$ has full column rank (there is no perfect multicollinearity);
5. $\text{var}(u_i | X_i) = \sigma_u^2$ (homoskedasticity); and
6. The conditional distribution of $u_i$ given $X_i$ is normal (normal errors).

In Equation (18.3), the first regressor is the "constant" regressor that always equals 1, and its coefficient is the intercept. Thus the intercept does not appear separately in Equation (18.3); rather, it is the first element of the coefficient vector $\boldsymbol{\beta}$.

Stacking all $n$ observations in Equation (18.3) yields the multiple regression model in matrix form:

$$Y = X\boldsymbol{\beta} + U. \tag{18.5}$$

## The Extended Least Squares Assumptions

The extended least squares assumptions for the multiple regressor model are the four least squares assumptions for the multiple regression model in Key Concept 6.4, plus the two additional assumptions of homoskedasticity and normally distributed errors. The assumption of homoskedasticity is used when we study the efficiency of the OLS estimator, and the assumption of normality is used when we study the exact sampling distribution of the OLS estimator and test statistics.

The extended least squares assumptions are summarized in Key Concept 18.1.

Except for notational differences, the first three assumptions in Key Concept 18.1 are identical to the first three assumptions in Key Concept 6.4.

The fourth assumption in Key Concepts 6.4 and 18.1 might appear different, but in fact they are the same: They are simply different ways of saying that there cannot be perfect multicollinearity. Recall that perfect multicollinearity arises when one regressor can be written as a perfect linear combination of the others. In the matrix notation of Equation (18.2), perfect multicollinearity means that one column of $X$ is a perfect linear combination of the other columns of $X$, but if this is true, then $X$ does not have full column rank. Thus saying that $X$ has rank $k + 1$, that is, rank equal to the number of columns of $X$, is just another way to say that the regressors are not perfectly multicollinear.

The fifth least squares assumption in Key Concept 18.1 is that the error term is conditionally homoskedastic, and the sixth assumption is that the conditional distribution of $u_i$, given $X_i$, is normal. These two assumptions are the same as the final two assumptions in Key Concept 17.1, except that they are now stated for multiple regressors.

*Implications for the mean vector and covariance matrix of U.*  The least squares assumptions in Key Concept 18.1 imply simple expressions for the mean vector and covariance matrix of the conditional distribution of $U$ given the matrix of regressors $X$. (The mean vector and covariance matrix of a vector of random

variables are defined in Appendix 18.2.) Specifically, the first and second assumptions in Key Concept 18.1 imply that $E(u_i|X) = E(u_i|X_i) = 0$ and that $\text{cov}(u_i, u_j|X) = E(\boldsymbol{u_i u_j}|X) = E(u_i u_j|X_i, X_j) = E(u_i|X_i)E(u_j|X_j) = 0$ for $i \neq j$ (Exercise 17.7). The first, second, and fifth assumptions imply that $E(u_i^2|X) = E(u_i^2|X_i) = \sigma_u^2$. Combining these results, we have that

$$\text{under Assumptions \#1 and \#2, } E(\boldsymbol{U}|\boldsymbol{X}) = \boldsymbol{0}_n, \text{ and} \tag{18.6}$$

$$\text{under Assumptions \#1, \#2, and \#5, } E(\boldsymbol{U}\boldsymbol{U}'|\boldsymbol{X}) = \sigma_u^2\boldsymbol{I}_n, \tag{18.7}$$

where $\boldsymbol{0}_n$ is the $n$-dimensional vector of zeros and $\boldsymbol{I}_n$ is the $n \times n$ identity matrix.

Similarly, the first, second, fifth, and sixth assumptions in Key Concept 18.1 imply that the conditional distribution of the $n$-dimensional random vector $\boldsymbol{U}$, conditional on $\boldsymbol{X}$, is the multivariate normal distribution (defined in Appendix 18.2). That is,

$$\begin{aligned}&\text{under Assumptions \#1, \#2, \#5, and \#6, the}\\&\text{conditional distribution of } \boldsymbol{U} \text{ given } \boldsymbol{X} \text{ is } N(\boldsymbol{0}_n, \sigma_u^2\boldsymbol{I}_n).\end{aligned} \tag{18.8}$$

## The OLS Estimator

The OLS estimator minimizes the sum of squared prediction mistakes, $\sum_{i=1}^{n}(Y_i - b_0 - b_1X_{1i} - \cdots - b_kX_{ki})^2$ [Equation (6.8)]. The formula for the OLS estimator is obtained by taking the derivative of the sum of squared prediction mistakes with respect to each element of the coefficient vector, setting these derivatives to zero, and solving for the estimator $\hat{\boldsymbol{\beta}}$.

The derivative of the sum of squared prediction mistakes with respect to the $j^{\text{th}}$ regression coefficient, $b_j$, is

$$\frac{\partial}{\partial b_j}\sum_{i=1}^{n}(Y_i - b_0 - b_1X_{1i} - \cdots - b_kX_{ki})^2$$

$$= -2\sum_{i=1}^{n}X_{ji}(Y_i - b_0 - b_1X_{1i} - \cdots - b_kX_{ki}) \tag{18.9}$$

for $j = 0, \ldots, k$, where, for $j = 0$, $X_{0i} = 1$ for all $i$. The derivative on the right-hand side of Equation (18.9) is the $j^{\text{th}}$ element of the $k + 1$ dimensional vector, $-2\boldsymbol{X}'(\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{b})$, where $\boldsymbol{b}$ is the $k + 1$ dimensional vector consisting of $b_0, \ldots, b_k$. There are $k + 1$ such derivatives, each corresponding to an element of $\boldsymbol{b}$. Combined, these yield the system of $k + 1$ equations that, when set to zero, constitute

the first order conditions for the OLS estimator $\hat{\boldsymbol{\beta}}$. That is, $\hat{\boldsymbol{\beta}}$ solves the system of $k + 1$ equations

$$\boldsymbol{X}'(\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}}) = \boldsymbol{0}_{k+1}, \tag{18.10}$$

or, equivalently, $\boldsymbol{X}'\boldsymbol{Y} = \boldsymbol{X}'\boldsymbol{X}\hat{\boldsymbol{\beta}}$.

Solving the system of equations (18.10) yields the OLS estimator $\hat{\boldsymbol{\beta}}$ in matrix form:

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{Y}, \tag{18.11}$$

where $(\boldsymbol{X}'\boldsymbol{X})^{-1}$ is the inverse of the matrix $\boldsymbol{X}'\boldsymbol{X}$.

***The role of "no perfect multicollinearity."*** The fourth least squares assumption in Key Concept 18.1 states that $\boldsymbol{X}$ has full column rank. In turn, this implies that the matrix $\boldsymbol{X}'\boldsymbol{X}$ has full rank, that is, $\boldsymbol{X}'\boldsymbol{X}$ is nonsingular. Because $\boldsymbol{X}'\boldsymbol{X}$ is nonsingular, it is invertible. Thus the assumption that there is no perfect multicollinearity ensures that $(\boldsymbol{X}'\boldsymbol{X})^{-1}$ exists, so Equation (18.10) has a unique solution and the formula in Equation (18.11) for the OLS estimator can actually be computed. Said differently, if $\boldsymbol{X}$ does *not* have full column rank, there is not a unique solution to Equation (18.10) and $\boldsymbol{X}'\boldsymbol{X}$ is singular. Therefore, $(\boldsymbol{X}'\boldsymbol{X})^{-1}$ cannot be computed and thus $\hat{\boldsymbol{\beta}}$ cannot be computed from Equation (18.11).

## 18.2 Asymptotic Distribution of the OLS Estimator and *t*-Statistic

If the sample size is large and the first four assumptions of Key Concept 18.1 are satisfied, then the OLS estimator has an asymptotic joint normal distribution, the heteroskedasticity-robust estimator of the covariance matrix is consistent, and the heteroskedasticity-robust OLS *t*-statistic has an asymptotic standard normal distribution. These results make use of the multivariate normal distribution (Appendix 18.2) and a multivariate extension of the central limit theorem.

### The Multivariate Central Limit Theorem

The central limit theorem of Key Concept 2.7 applies to a one-dimensional random variable. To derive the *joint* asymptotic distribution of the elements of $\hat{\boldsymbol{\beta}}$, we need a multivariate central limit theorem that applies to vector-valued random variables.

## The Multivariate Central Limit Theorem

Suppose that $W_1, \ldots, W_n$ are i.i.d. $m$-dimensional random variables with mean vector $E(W_i) = \mu_W$ and covariance matrix $E[(W_i - \mu_W)(W_i - \mu_W)'] = \Sigma_W$, where $\Sigma_W$ is positive definite and finite. Let $\overline{W} = \frac{1}{n}\sum_{i=1}^{n}W_i$. Then $\sqrt{n}(\overline{W} - \mu_W) \xrightarrow{d} N(\mathbf{0}_m, \Sigma_W)$.

The multivariate central limit theorem extends the univariate central limit theorem to averages of observations on a vector-valued random variable, $W$, where $W$ is $m$-dimensional. The difference between the central limit theorems for a scalar as opposed to a vector-valued random variable is the conditions on the variances. In the scalar case in Key Concept 2.7, the requirement is that the variance is both nonzero and finite. In the vector case, the requirement is that the covariance matrix is both positive definite and finite. If the vector-valued random variable $W$ has a finite positive definite covariance matrix, then $0 < \text{var}(c'W) < \infty$ for all nonzero $m$-dimensional vectors $c$ (Exercise 18.3).

The multivariate central limit theorem that we will use is stated in Key Concept 18.2.

### Asymptotic Normality of $\hat{\beta}$

In large samples, the OLS estimator has the multivariate normal asymptotic distribution

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}_{k+1}, \Sigma_{\sqrt{n}(\hat{\beta} - \beta)}), \text{ where } \Sigma_{\sqrt{n}(\hat{\beta} - \beta)} = Q_X^{-1}\Sigma_V Q_X^{-1}, \quad (18.12)$$

where $Q_X$ is the $(k + 1) \times (k + 1)$-dimensional matrix of second moments of the regressors, that is, $Q_X = E(X_iX_i')$, and $\Sigma_V$ is the $(k + 1) \times (k + 1)$-dimensional covariance matrix of $V_i = X_iu_i$, that is, $\Sigma_V = E(V_iV_i')$. Note that the second least squares assumption in Key Concept 18.1 implies that $V_i, i = 1, \ldots, n$, are i.i.d. Written in terms of $\hat{\beta}$ rather than $\sqrt{n}(\hat{\beta} - \beta)$, the normal approximation in Equation (18.12) is

$$\hat{\beta}, \text{ in large samples, is approximately distributed } N(\beta, \Sigma_{\hat{\beta}})$$
$$\text{where } \Sigma_{\hat{\beta}} = \Sigma_{\sqrt{n}(\hat{\beta} - \beta)}/n = Q_X^{-1}\Sigma_V Q_X^{-1}/n. \quad (18.13)$$

The covariance matrix $\Sigma_{\hat{\beta}}$ in Equation (18.13) is the covariance matrix of the approximate normal distribution of $\hat{\beta}$, whereas $\Sigma_{\sqrt{n}(\hat{\beta}-\beta)}$ in Equation (18.12) is the covariance matrix of the asymptotic normal distribution of $\sqrt{n}(\hat{\beta}-\beta)$. These two covariance matrices differ by a factor of $n$, depending on whether the OLS estimator is scaled by $\sqrt{n}$.

***Derivation of Equation (18.12).***  To derive Equation (18.12), first use Equations (18.4) and (18.11) to write $\hat{\beta} = (X'X)^{-1}X'Y = (X'X)^{-1}X'(X\beta + U)$ so that

$$\hat{\beta} = \beta + (X'X)^{-1}X'U. \tag{18.14}$$

Thus $\hat{\beta} - \beta = (X'X)^{-1}X'U$, so

$$\sqrt{n}(\hat{\beta} - \beta) = \left(\frac{X'X}{n}\right)^{-1}\left(\frac{X'U}{\sqrt{n}}\right). \tag{18.15}$$

The derivation of Equation (18.12) involves arguing first that the "denominator" matrix in Equation (18.15), $X'X/n$, is consistent for $Q_X$ and second that the "numerator" matrix, $X'U/\sqrt{n}$, obeys the multivariate central limit theorem in Key Concept 18.2. The details are given in Appendix 18.3.

## Heteroskedasticity-Robust Standard Errors

The heteroskedasticity-robust estimator of $\Sigma_{\sqrt{n}(\hat{\beta}-\beta)}$ is obtained by replacing the population moments in its definition [Equation (18.12)] by sample moments. Accordingly, the heteroskedasticity-robust estimator of the covariance matrix of $\sqrt{n}(\hat{\beta}-\beta)$ is

$$\hat{\Sigma}_{\sqrt{n}(\hat{\beta}-\beta)} = \left(\frac{X'X}{n}\right)^{-1}\hat{\Sigma}_{\hat{V}}\left(\frac{X'X}{n}\right)^{-1}, \text{ where } \hat{\Sigma}_{\hat{V}} = \frac{1}{n-k-1}\sum_{i=1}^{n}X_iX_i'\hat{u}_i^2, \tag{18.16}$$

The estimator $\hat{\Sigma}_{\hat{V}}$ incorporates the same degrees-of-freedom adjustment that is in the *SER* for the multiple regression model (Section 6.4) to adjust for potential downward bias because of estimation of $k + 1$ regression coefficients.

The proof that $\hat{\Sigma}_{\sqrt{n}(\hat{\beta}-\beta)} \xrightarrow{p} \Sigma_{\sqrt{n}(\hat{\beta}-\beta)}$ is conceptually similar to the proof, presented in Section 17.3, of the consistency of heteroskedasticity-robust standard errors for the single-regressor model.

***Heteroskedasticity-robust standard errors.***  The heteroskedasticity-robust estimator of the covariance matrix of $\hat{\beta}$, $\Sigma_{\hat{\beta}}$ is

$$\hat{\Sigma}_{\hat{\beta}} = n^{-1}\hat{\Sigma}_{\sqrt{n}(\hat{\beta}-\beta)}. \tag{18.17}$$

The heteroskedasticity-robust standard error for the $j^{th}$ regression coefficient is the square root of the $j^{th}$ diagonal element of $\hat{\Sigma}_{\hat{\beta}}$. That is, the heteroskedasticity-robust standard error of the $j^{th}$ coefficient is

$$SE(\hat{\beta}_j) = \sqrt{(\hat{\Sigma}_{\hat{\beta}})_{jj}}, \tag{18.18}$$

where $(\hat{\Sigma}_{\hat{\beta}})_{jj}$ is the $(j, j)$ element of $\hat{\Sigma}_{\hat{\beta}}$.

## Confidence Intervals for Predicted Effects

Section 8.1 describes two methods for computing the standard error of predicted effects that involve changes in two or more regressors. There are compact matrix expressions for these standard errors and thus for confidence intervals for predicted effects.

Consider a change in the value of the regressors for the $i^{th}$ observation from some initial value, say $X_{i,0}$, to some new value, $X_{i,0} + d$, so that the change in $X_i$ is $\Delta X_i = d$, where $d$ is a $k + 1$ dimensional vector. This change in $X$ can involve multiple regressors (that is, multiple elements of $X_i$). For example, if two of the regressors are the value of an independent variable and its square, then $d$ is the difference between the subsequent and initial values of these two variables.

The expected effect of this change in $X_i$ is $d'\beta$, and the estimator of this effect is $d'\hat{\beta}$. Because linear combinations of normally distributed random variables are themselves normally distributed, $\sqrt{n}(d'\hat{\beta} - d'\beta) = d'\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, d'\Sigma_{\sqrt{n}(\hat{\beta} - \beta)}d)$. Thus the standard error of this predicted effect is $(d'\hat{\Sigma}_{\hat{\beta}}d)^{1/2}$. A 95% confidence interval for this predicted effect is

$$d'\hat{\beta} \pm 1.96\sqrt{d'\hat{\Sigma}_{\hat{\beta}}d}. \tag{18.19}$$

## Asymptotic Distribution of the $t$-Statistic

The $t$-statistic testing the null hypothesis that $\beta_j = \beta_{j,0}$, constructed using the heteroskedasticity-robust standard error in Equation (18.18), is given in Key Concept 7.1. The argument that this $t$-statistic has an asymptotic standard normal distribution parallels the argument given in Section 17.3 for the single-regressor model.

# 18.3  Tests of Joint Hypotheses

Section 7.2 considers tests of joint hypotheses that involve multiple restrictions, where each restriction involves a single coefficient, and Section 7.3 considers tests of a single restriction involving two or more coefficients. The matrix setup of

Section 18.1 permits a unified representation of these two types of hypotheses as linear restrictions on the coefficient vector, where each restriction can involve multiple coefficients. Under the first four least squares assumptions in Key Concept 18.1, the heteroskedasticity-robust OLS $F$-statistic testing these hypotheses has an $F_{q,\infty}$ asymptotic distribution under the null hypothesis.

## Joint Hypotheses in Matrix Notation

Consider a joint hypothesis that is linear in the coefficients and imposes $q$ restrictions, where $q \leq k + 1$. Each of these $q$ restrictions can involve one or more of the regression coefficients. This joint null hypothesis can be written in matrix notation as

$$R\beta = r, \tag{18.20}$$

where $R$ is a $q \times (k + 1)$ nonrandom matrix with full row rank and $r$ is a nonrandom $q \times 1$ vector. The number of rows of $R$ is $q$, which is the number of restrictions being imposed under the null hypothesis.

The null hypothesis in Equation (18.20) subsumes all the null hypotheses considered in Sections 7.2 and 7.3. For example, a joint hypothesis of the type considered in Section 7.2 is that $\beta_0 = 0, \beta_1 = 0, \ldots, \beta_{q-1} = 0$. To write this joint hypothesis in the form of Equation (18.20), set $R = [I_q \quad 0_{q \times (k+1-q)}]$ and $r = 0_q$.

The formulation in Equation (18.20) also captures the restrictions of Section 7.3 involving multiple regression coefficients. For example, if $k = 2$, then the hypothesis that $\beta_1 + \beta_2 = 1$ can be written in the form of Equation (18.20) by setting $R = [0 \quad 1 \quad 1]$, $r = 1$, and $q = 1$.

## Asymptotic Distribution of the $F$-Statistic

The heteroskedasticity-robust $F$-statistic testing the joint hypothesis in Equation (18.20) is

$$F = (R\hat{\beta} - r)' [R\hat{\Sigma}_{\hat{\beta}} R']^{-1} (R\hat{\beta} - r)/q. \tag{18.21}$$

If the first four assumptions in Key Concept 18.1 hold, then under the null hypothesis

$$F \xrightarrow{d} F_{q,\infty}. \tag{18.22}$$

This result follows by combining the asymptotic normality of $\hat{\boldsymbol{\beta}}$ with the consistency of the heteroskedasticity-robust estimator $\hat{\boldsymbol{\Sigma}}_{\sqrt{n}(\hat{\boldsymbol{\beta}}-\boldsymbol{\beta})}$ of the covariance matrix. Specifically, first note that Equation (18.12) and Equation (18.74) in Appendix 18.2 imply that, under the null hypothesis, $\sqrt{n}(\boldsymbol{R}\hat{\boldsymbol{\beta}} - \boldsymbol{r}) = \sqrt{n}\boldsymbol{R}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} N(\boldsymbol{0}, \boldsymbol{R}\boldsymbol{\Sigma}_{\sqrt{n}(\hat{\boldsymbol{\beta}}-\boldsymbol{\beta})}\boldsymbol{R}')$. It follows from Equation (18.77) that, under the null hypothesis, $(\boldsymbol{R}\hat{\boldsymbol{\beta}} - \boldsymbol{r})'[\boldsymbol{R}\boldsymbol{\Sigma}_{\hat{\boldsymbol{\beta}}}\boldsymbol{R}']^{-1}(\boldsymbol{R}\hat{\boldsymbol{\beta}} - \boldsymbol{r}) = [\sqrt{n}\boldsymbol{R}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})]'[\boldsymbol{R}\boldsymbol{\Sigma}_{\sqrt{n}(\hat{\boldsymbol{\beta}}-\boldsymbol{\beta})}\boldsymbol{R}']^{-1}[\sqrt{n}\boldsymbol{R}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})] \xrightarrow{d} \chi_q^2$. However, because $\hat{\boldsymbol{\Sigma}}_{\sqrt{n}(\hat{\boldsymbol{\beta}}-\boldsymbol{\beta})} \xrightarrow{p} \boldsymbol{\Sigma}_{\sqrt{n}(\hat{\boldsymbol{\beta}}-\boldsymbol{\beta})}$, it follows from Slutsky's theorem that $[\sqrt{n}\boldsymbol{R}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})]'[\boldsymbol{R}\hat{\boldsymbol{\Sigma}}_{\sqrt{n}(\hat{\boldsymbol{\beta}}-\boldsymbol{\beta})}\boldsymbol{R}']^{-1}[\sqrt{n}\boldsymbol{R}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})] \xrightarrow{d} \chi_q^2$. or, equivalently (because $\hat{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{\beta}}} = \hat{\boldsymbol{\Sigma}}_{\sqrt{n}(\hat{\boldsymbol{\beta}}-\boldsymbol{\beta})}/n$), that $F \xrightarrow{d} \chi_q^2/q$, which is in turn distributed $F_{q,\infty}$.

## Confidence Sets for Multiple Coefficients

As discussed in Section 7.4, an asymptotically valid confidence set for two or more elements of $\boldsymbol{\beta}$ can be constructed as the set of values that, when taken as the null hypothesis, are not rejected by the $F$-statistic. In principle, this set could be computed by repeatedly evaluating the $F$-statistic for many values of $\boldsymbol{\beta}$, but, as is the case with a confidence interval for a single coefficient, it is simpler to manipulate the formula for the test statistic to obtain an explicit formula for the confidence set.

Here is the procedure for constructing a confidence set for two or more of the elements of $\boldsymbol{\beta}$. Let $\boldsymbol{\delta}$ denote the $q$-dimensional vector consisting of the coefficients for which we wish to construct a confidence set. For example, if we are constructing a confidence set for the regression coefficients $\beta_1$ and $\beta_2$, then $q = 2$ and $\boldsymbol{\delta} = (\beta_1 \beta_2)'$. In general, we can write $\boldsymbol{\delta} = \boldsymbol{R}\boldsymbol{\beta}$, where the matrix $\boldsymbol{R}$ consists of zeros and ones [as discussed following Equation (18.20)]. The $F$-statistic testing the hypothesis that $\boldsymbol{\delta} = \boldsymbol{\delta}_0$ is $F = (\hat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0)'[\boldsymbol{R}\hat{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{\beta}}}\boldsymbol{R}']^{-1}(\hat{\boldsymbol{\delta}} - \boldsymbol{\delta}_0)/q$, where $\hat{\boldsymbol{\delta}} = \boldsymbol{R}\hat{\boldsymbol{\beta}}$. A 95% confidence set for $\boldsymbol{\delta}$ is the set of values $\boldsymbol{\delta}_0$ that are not rejected by the $F$-statistic. That is, when $\boldsymbol{\delta} = \boldsymbol{R}\boldsymbol{\beta}$, a 95% confidence set for $\boldsymbol{\delta}$ is

$$\{\boldsymbol{\delta}: (\hat{\boldsymbol{\delta}} - \boldsymbol{\delta})'[\boldsymbol{R}\hat{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{\beta}}}\boldsymbol{R}']^{-1}(\hat{\boldsymbol{\delta}} - \boldsymbol{\delta})/q \leq c\}, \tag{18.23}$$

where $c$ is the 95[th] percentile (the 5% critical value) of the $F_{q,\infty}$ distribution.

The set in Equation (18.23) consists of all the points contained inside the ellipse determined when the inequality in Equation (18.23) is an equality (this is an ellipsoid when $q > 2$). Thus the confidence set for $\delta$ can be computed by solving Equation (18.23) for the boundary ellipse.

## 18.4 Distribution of Regression Statistics with Normal Errors

The distributions presented in Sections 18.2 and 18.3, which were justified by appealing to the law of large numbers and the central limit theorem, apply when the sample size is large. If, however, the errors are homoskedastic and normally distributed, conditional on $X$, then the OLS estimator has a multivariate normal distribution in finite sample, conditional on $X$. In addition, the finite sample distribution of the square of the standard error of the regression is proportional to the chi-squared distribution with $n - k - 1$ degrees of freedom, the homoskedasticity-only OLS $t$-statistic has a Student $t$ distribution with $n - k - 1$ degrees of freedom, and the homoskedasticity-only $F$-statistic has an $F_{q,n-k-1}$ distribution. The arguments in this section employ some specialized matrix formulas for OLS regression statistics, which are presented first.

### Matrix Representations of OLS Regression Statistics

The OLS predicted values, residuals, and sum of squared residuals have compact matrix representations. These representations make use of two matrices, $P_X$ and $M_X$.

***The matrices $P_X$ and $M_X$.*** The algebra of OLS in the multivariate model relies on the two symmetric $n \times n$ matrices, $P_X$ and $M_X$:

$$P_X = X(X'X)^{-1}X' \text{ and} \qquad (18.24)$$

$$M_X = I_n - P_X. \qquad (18.25)$$

A matrix $C$ is idempotent if $C$ is square and $CC = C$ (see Appendix 18.1). Because $P_X = P_X P_X$ and $M_X = M_X M_X$ (Exercise 18.5), and because $P_X$ and $M_X$ are symmetric, $P_X$ and $M_X$ are symmetric idempotent matrices.

The matrices $P_X$ and $M_X$ have some additional useful properties (Exercise 18.5), which follow directly from the definitions in Equations (18.24) and (18.25):

$$P_X X = X \text{ and } M_X X = 0_{n \times (k+1)};$$
$$\text{rank}(P_X) = k + 1 \text{ and rank}(M_X) = n - k - 1, \qquad (18.26)$$

where $\text{rank}(P_X)$ is the rank of $P_X$.

The matrices $P_X$ and $M_X$ can be used to decompose an $n$-dimensional vector $Z$ into two parts: a part that is spanned by the columns of $X$ and a part orthogonal to the columns of $X$. In other words, $P_XZ$ is the projection of $Z$ onto the space spanned by the columns of $X$, $M_XZ$ is the part of $Z$ orthogonal to the columns of $X$, and $Z = P_XZ + M_XZ$.

*OLS predicted values and residuals.* The matrices $P_X$ and $M_X$ provide some simple expressions for OLS predicted values and residuals. The OLS predicted values, $\hat{Y} = X\hat{\beta}$, and the OLS residuals, $\hat{U} = Y - \hat{Y}$, can be expressed as follows (Exercise 18.5):

$$\hat{Y} = P_XY \text{ and} \tag{18.27}$$

$$\hat{U} = M_XY = M_XU. \tag{18.28}$$

The expressions in Equations (18.27) and (18.28) provide a simple proof that the OLS residuals and predicted values are orthogonal, that is, Equation (4.37) holds: $\hat{Y}'\hat{U} = Y'P_X'M_XY = 0$, where the second equality follows from $P_X'M_X = \mathbf{0}_{n \times n}$, which in turn follows from $M_XX = \mathbf{0}_{n \times (k+1)}$ in Equation (18.26).

*The standard error of the regression.* The *SER*, defined in Section 4.3, is $s_{\hat{u}}$, where

$$s_{\hat{u}}^2 = \frac{1}{n-k-1}\sum_{i=1}^{n}\hat{u}_i^2 = \frac{1}{n-k-1}\hat{U}'\hat{U} = \frac{1}{n-k-1}U'M_XU, \tag{18.29}$$

where the final equality follows because $\hat{U}'\hat{U} = (M_XU)'(M_XU) = U'M_XM_XU = U'M_XU$ (because $M_X$ is symmetric and idempotent).

## Distribution of $\hat{\beta}$ with Normal Errors

Because $\hat{\beta} = \beta + (X'X)^{-1}X'U$ [Equation (18.14)] and because the distribution of $U$ conditional on $X$ is, by assumption, $N(\mathbf{0}_n, \sigma_u^2I_n)$ [Equation (18.8)], the conditional distribution of $\hat{\beta}$ given $X$ is multivariate normal with mean $\beta$. The covariance matrix of $\hat{\beta}$, conditional on $X$, is $\Sigma_{\hat{\beta}|X} = E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'|X] = E[(X'X)^{-1}X'UU'X(X'X)^{-1}|X] = (X'X)^{-1}X'(\sigma_u^2I_n)X(X'X)^{-1} = \sigma_u^2(X'X)^{-1}$.

Accordingly, under all six assumptions in Key Concept 18.1, the finite-sample conditional distribution of $\hat{\boldsymbol{\beta}}$ given $X$ is

$$\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \boldsymbol{\Sigma}_{\hat{\beta}|X}), \text{ where } \boldsymbol{\Sigma}_{\hat{\beta}|X} = \sigma_u^2(X'X)^{-1}. \tag{18.30}$$

## Distribution of $s_{\hat{u}}^2$

If all six assumptions in Key Concept 18.1 hold, then $s_{\hat{u}}^2$ has an exact sampling distribution that is proportional to a chi-squared distribution with $n - k - 1$ degrees of freedom:

$$s_{\hat{u}}^2 \sim \frac{\sigma_u^2}{n - k - 1} \times \chi_{n-k-1}^2 \tag{18.31}$$

The proof of Equation (18.31) starts with Equation (18.29). Because $U$ is normally distributed conditional on $X$ and because $M_X$ is a symmetric idempotent matrix, the quadratic form $U'M_XU/\sigma_u^2$ has an exact chi-squared distribution with degrees of freedom equal to the rank of $M_X$ [Equation (18.78) in Appendix 18.2]. From Equation (18.26), the rank of $M_X$ is $n - k - 1$. Thus $U'M_XU/\sigma_u^2$ has an exact $\chi_{n-k-1}^2$ distribution, from which Equation (18.31) follows.

The degrees-of-freedom adjustment ensures that $s_{\hat{u}}^2$ is unbiased. The expectation of a random variable with a $\chi_{n-k-1}^2$ distribution is $n - k - 1$; thus $E(U'M_XU) = (n - k - 1)\sigma_u^2$, so $E(s_{\hat{u}}^2) = \sigma_u^2$.

## Homoskedasticity-Only Standard Errors

The homoskedasticity-only estimator $\widetilde{\boldsymbol{\Sigma}}_{\hat{\beta}}$ of the covariance matrix of $\hat{\boldsymbol{\beta}}$, conditional on $X$, is obtained by substituting the sample variance $s_{\hat{u}}^2$ for the population variance $\sigma_u^2$ in the expression for $\boldsymbol{\Sigma}_{\hat{\beta}|X}$ in Equation (18.30). Accordingly,

$$\widetilde{\boldsymbol{\Sigma}}_{\hat{\beta}} = s_{\hat{u}}^2(X'X)^{-1} \quad \text{(homoskedasticity-only).} \tag{18.32}$$

The estimator of the variance of the normal conditional distribution of $\hat{\beta}_j$, given $X$, is the $(j, j)$ element of $\widetilde{\boldsymbol{\Sigma}}_{\hat{\beta}}$. Thus the homoskedasticity-only standard error of $\hat{\beta}_j$ is the square root of the $j^{\text{th}}$ diagonal element of $\widetilde{\boldsymbol{\Sigma}}_{\hat{\beta}}$. That is, the homoskedasticity-only standard error of $\hat{\beta}_j$ is

$$\widetilde{SE}(\hat{\beta}_j) = \sqrt{(\widetilde{\boldsymbol{\Sigma}}_{\hat{\beta}})_{jj}} \quad \text{(homoskedasticity-only).} \tag{18.33}$$

## Distribution of the *t*-Statistic

Let $\tilde{t}$ be the *t*-statistic testing the hypothesis $\beta_j = \beta_{j,0}$, constructed using the homoskedasticity-only standard error; that is, let

$$\tilde{t} = \frac{\hat{\beta}_j - \beta_{j,0}}{\sqrt{(\widetilde{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{\beta}}})_{jj}}}. \tag{18.34}$$

Under all six of the extended least squares assumptions in Key Concept 18.1, the exact sampling distribution of $\tilde{t}$ is the Student *t* distribution with $n - k - 1$ degrees of freedom; that is,

$$\tilde{t} \sim t_{n-k-1}. \tag{18.35}$$

The proof of Equation (18.35) is given in Appendix 18.4.

## Distribution of the *F*-Statistic

If all six least squares assumptions in Key Concept 18.1 hold, then the *F*-statistic testing the hypothesis in Equation (18.20), constructed using the homoskedasticity-only estimator of the covariance matrix, has an exact $F_{q,\,n-k-1}$ distribution under the null hypothesis.

***The homoskedasticity-only F-statistic.*** The homoskedasticity-only *F*-statistic is similar to the heteroskedasticity-robust *F*-statistic in Equation (18.21), except that the homoskedasticity-only estimator $\widetilde{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{\beta}}}$ is used instead of the heteroskedasticity-robust estimator $\widetilde{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{\beta}}}$. Substituting the expression $\widetilde{\boldsymbol{\Sigma}}_{\hat{\boldsymbol{\beta}}} = s_{\hat{u}}^2 (\boldsymbol{X}'\boldsymbol{X})^{-1}$ into the expression for the *F*-statistic in Equation (18.21) yields the homoskedasticity-only *F*-statistic testing the null hypothesis in Equation (18.20):

$$\tilde{F} = \frac{(\boldsymbol{R}\hat{\boldsymbol{\beta}} - \boldsymbol{r})' [\boldsymbol{R}(\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{R}']^{-1}(\boldsymbol{R}\hat{\boldsymbol{\beta}} - \boldsymbol{r})/q}{s_{\hat{u}}^2}. \tag{18.36}$$

If all six assumptions in Key Concept 18.1 hold, then under the null hypothesis

$$\tilde{F} \sim F_{q,n-k-1}. \tag{18.37}$$

The proof of Equation (18.37) is given in Appendix 18.4.

The *F*-statistic in Equation (18.36) is called the Wald version of the *F*-statistic (named after the statistician Abraham Wald). Although the formula for the homoskedastic-only *F*-statistic given in Equation (7.13) appears quite different from the formula for the Wald statistic in Equation (18.36), the homoskedastic-only *F*-statistic and the Wald *F*-statistic are two versions of the same statistic. That is, the two expressions are equivalent, a result shown in Exercise 18.13.

## 18.5 Efficiency of the OLS Estimator with Homoskedastic Errors

Under the Gauss–Markov conditions for multiple regression, the OLS estimator of $\boldsymbol{\beta}$ is efficient among all linear conditionally unbiased estimators; that is, the OLS estimator is BLUE.

### The Gauss–Markov Conditions for Multiple Regression

The **Gauss–Markov conditions for multiple regression** are

$$(i)\ E(\boldsymbol{U}|\boldsymbol{X}) = \boldsymbol{0}_n,$$

$$(ii)\ E(\boldsymbol{U}\boldsymbol{U}'|\boldsymbol{X}) = \sigma_u^2 \boldsymbol{I}_n,\ \text{and}$$

$$(iii)\ \boldsymbol{X}\ \text{has full column rank.} \tag{18.38}$$

The Gauss–Markov conditions for multiple regression in turn are implied by the first five assumptions in Key Concept 18.1 [see Equations (18.6) and (18.7)]. The conditions in Equation (18.38) generalize the Gauss–Markov conditions for a single regressor model to multiple regression. [By using matrix notation, the second and third Gauss–Markov conditions in Equation (5.31) are collected into the single condition (ii) in Equation (18.38).]

### Linear Conditionally Unbiased Estimators

We start by describing the class of linear unbiased estimators and by showing that OLS is in that class.

***The class of linear conditionally unbiased estimators.*** An estimator of $\boldsymbol{\beta}$ is said to be linear if it is a linear function of $Y_1, \ldots, Y_n$. Accordingly, the estimator $\widetilde{\boldsymbol{\beta}}$ is linear in $\boldsymbol{Y}$ if it can be written in the form

$$\widetilde{\boldsymbol{\beta}} = \boldsymbol{A}'\boldsymbol{Y}, \tag{18.39}$$

where $A$ is an $n \times (k + 1)$ dimensional matrix of weights that may depend on $X$ and on nonrandom constants, but not on $Y$.

An estimator is conditionally unbiased if the mean of its conditional sampling distribution, given $X$, is $\boldsymbol{\beta}$. That is, $\widetilde{\boldsymbol{\beta}}$ is conditionally unbiased if $E(\widetilde{\boldsymbol{\beta}}|X) = \boldsymbol{\beta}$.

***The OLS estimator is linear and conditionally unbiased.***   Comparison of Equations (18.11) and (18.39) shows that the OLS estimator is linear in $Y$; specifically, $\hat{\boldsymbol{\beta}} = \hat{A}'Y$, where $\hat{A} = X(X'X)^{-1}$. To show that $\hat{\boldsymbol{\beta}}$ is conditionally unbiased, recall from Equation (18.14) that $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + (X'X)^{-1}X'U$. Taking the conditional expectation of both sides of this expression yields, $E(\hat{\boldsymbol{\beta}}|X) = \boldsymbol{\beta} + E[(X'X)^{-1}X'U|X] = \boldsymbol{\beta} + (X'X)^{-1}X'E(U|X) = \boldsymbol{\beta}$, where the final equality follows because $E(U|X) = 0$ by the first Gauss–Markov condition.

## The Gauss–Markov Theorem for Multiple Regression

The **Gauss–Markov theorem for multiple regression** provides conditions under which the OLS estimator is efficient among the class of linear conditionally unbiased estimators. A subtle point arises, however, because $\hat{\boldsymbol{\beta}}$ is a vector and its "variance" is a covariance matrix. When the "variance" of an estimator is a matrix, just what does it mean to say that one estimator has a smaller variance than another?

The Gauss–Markov theorem handles this problem by comparing the variance of a candidate estimator of a *linear combination* of the elements of $\boldsymbol{\beta}$ to the variance of the corresponding linear combination of $\hat{\boldsymbol{\beta}}$. Specifically, let $c$ be a $k + 1$ dimensional vector and consider the problem of estimating the linear combination $c'\boldsymbol{\beta}$ using the candidate estimator $c'\widetilde{\boldsymbol{\beta}}$ (where $\widetilde{\boldsymbol{\beta}}$ is a linear conditionally unbiased estimator) on the one hand and $c'\hat{\boldsymbol{\beta}}$ on the other hand. Because $c'\widetilde{\boldsymbol{\beta}}$ and $c'\hat{\boldsymbol{\beta}}$ are both scalars and are both linear conditionally unbiased estimators of $c'\boldsymbol{\beta}$, it now makes sense to compare their variances.

The Gauss–Markov theorem for multiple regression says that the OLS estimator of $c'\boldsymbol{\beta}$ is efficient; that is, the OLS estimator $c'\hat{\boldsymbol{\beta}}$ has the smallest conditional variance of all linear conditionally unbiased estimators $c'\widetilde{\boldsymbol{\beta}}$. Remarkably, this is true no matter what the linear combination is. It is in this sense that the OLS estimator is BLUE in multiple regression.

The Gauss–Markov theorem is stated in Key Concept 18.3 and proven in Appendix 18.5.

### Gauss–Markov Theorem for Multiple Regression

Suppose that the Gauss–Markov conditions for multiple regression in Equation (18.38) hold. Then the OLS estimator $\hat{\boldsymbol{\beta}}$ is BLUE. That is, let $\tilde{\boldsymbol{\beta}}$ be a linear conditionally unbiased estimator of $\boldsymbol{\beta}$ and let $\boldsymbol{c}$ be a nonrandom $k + 1$ dimensional vector. Then $\text{var}(\boldsymbol{c}'\hat{\boldsymbol{\beta}}|\boldsymbol{X}) \leq \text{var}(\boldsymbol{c}'\tilde{\boldsymbol{\beta}}|\boldsymbol{X})$ for every nonzero vector $\boldsymbol{c}$, where the inequality holds with equality for all $\boldsymbol{c}$ only if $\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}$.

## 18.6  Generalized Least Squares[1]

The assumption of i.i.d. sampling fits many applications. For example, suppose that $Y_i$ and $X_i$ correspond to information about individuals, such as their earnings, education, and personal characteristics, where the individuals are selected from a population by simple random sampling. In this case, because of the simple random sampling scheme, $(X_i, Y_i)$ are necessarily i.i.d. Because $(X_i, Y_i)$ and $(X_j, Y_j)$ are independently distributed for $i \neq j$, $u_i$ and $u_j$ are independently distributed for $i \neq j$. This in turn implies that $u_i$ and $u_j$ are uncorrelated for $i \neq j$. In the context of the Gauss–Markov assumptions, the assumption that $E(\boldsymbol{U}\boldsymbol{U}'|\boldsymbol{X})$ is diagonal therefore is appropriate if the data are collected in a way that makes the observations independently distributed.

Some sampling schemes encountered in econometrics do not, however, result in independent observations and instead can lead to error terms $u_i$ that are correlated from one observation to the next. The leading example is when the data are sampled over time for the same entity, that is, when the data are time series data. As discussed in Section 15.3, in regressions involving time series data, many omitted factors are correlated from one period to the next, and this can result in regression error terms (which represent those omitted factors) that are correlated from one period of observation to the next. In other words, the error term in one period will not, in general, be distributed independently of the error term in the

---

[1]The GLS estimator was introduced in Section 15.5 in the context of distributed lag time series regression. This presentation here is a self-contained mathematical treatment of GLS that can be read independently of Section 15.5, but reading that section first will help to make these ideas more concrete.

next period. Instead, the error term in one period could be correlated with the error term in the next period.

The presence of correlated error terms creates two problems for inference based on OLS. First, *neither* the heteroskedasticity-robust nor the homoskedasticity-only standard errors produced by OLS provide a valid basis for inference. The solution to this problem is to use standard errors that are robust to both heteroskedasticity and correlation of the error terms across observations. This topic—heteroskedasticity- and autocorrelation-consistent (HAC) covariance matrix estimation—is the subject of Section 15.4 and we do not pursue it further here.

Second, if the error term is correlated across observations, then $E(UU'|X)$ is not diagonal, the second Gauss–Markov condition in Equation (18.38) does not hold, and OLS is not BLUE. In this section we study an estimator, **generalized least squares (GLS)**, that is BLUE (at least asymptotically) when the conditional covariance matrix of the errors is no longer proportional to the identity matrix. A special case of GLS is weighted least squares, discussed in Section 17.5, in which the conditional covariance matrix is diagonal and the $i$th diagonal element is a function of $X_i$. Like WLS, GLS transforms the regression model so that the errors of the transformed model satisfy the Gauss–Markov conditions. The GLS estimator is the OLS estimator of the coefficients in the transformed model.

## The GLS Assumptions

There are four assumptions under which GLS is valid. The first GLS assumption is that $u_i$ has a mean of zero, conditional on $X_1, \ldots, X_n$; that is,

$$E(U|X) = \mathbf{0}_n. \tag{18.40}$$

This assumption is implied by the first two least squares assumptions in Key Concept 18.1; that is, if $E(u_i|X_i) = 0$ and $(X_i, Y_i)$, $i = 1, \ldots, n$, are i.i.d., then $E(U|X) = \mathbf{0}_n$. In GLS, however, we will not want to maintain the i.i.d. assumption; after all, one purpose of GLS is to handle errors that are correlated across observations. We discuss the significance of the assumption in Equation (18.40) after introducing the GLS estimator.

The second GLS assumption is that the conditional covariance matrix of $U$ given $X$ is some function of $X$:

$$E(UU'|X) = \mathbf{\Omega}(X), \tag{18.41}$$

where $\mathbf{\Omega}(X)$ is an $n \times n$ positive definite matrix-valued function of $X$.

## The GLS Assumptions

In the linear regression model $Y = X\beta + U$, the GLS assumptions are

1. $E(U|X) = \mathbf{0}_n$;
2. $E(UU'|X) = \Omega(X)$, where $\Omega(X)$ is an $n \times n$ positive definite matrix that can depend on $X$;
3. $X_i$ and $u_i$ satisfy suitable moment conditions; and
4. $X$ has full column rank (there is no perfect multicollinearity).

There are two main applications of GLS that are covered by this assumption. The first is independent sampling with heteroskedastic errors, in which case $\Omega(X)$ is a diagonal matrix with diagonal element $\lambda h(X_i)$, where $\lambda$ is a constant and $h$ is a function. In this case, discussed in Section 17.5, GLS is WLS.

The second application is to homoskedastic errors that are serially correlated. In practice, in this case a model is developed for the serial correlation. For example, one model is that the error term is correlated with only its neighbor, so $\text{corr}(u_i, u_{i-1}) = \rho \neq 0$ but $\text{corr}(u_i, u_j) = 0$ if $|i - j| \geq 2$. In this case, $\Omega(X)$ has $\sigma_u^2$ as its diagonal element, $\rho\sigma_u^2$ in the first off-diagonal, and zeros elsewhere. Thus $\Omega(X)$ does not depend on $X$, $\Omega_{ii} = \sigma_u^2$, $\Omega_{ij} = \rho\sigma_u^2$ for $|i - j| = 1$, and $\Omega_{ij} = 0$ for $|i - j| > 1$. Other models for serial correlation, including the first order autoregressive model, are discussed further in the context of GLS in Section 15.5 (also see Exercise 18.8).

One assumption that has appeared on all previous lists of least squares assumptions for cross-sectional data is that $X_i$ and $u_i$ have nonzero, finite fourth moments. In the case of GLS, the specific moment assumptions needed to prove asymptotic results depend on the nature of the function $\Omega(X)$, whether $\Omega(X)$ is known or estimated, and the statistic under consideration (the GLS estimator, $t$-statistic, etc.). Because the assumptions are case- and model-specific, we do not present specific moment assumptions here, and the discussion of the large-sample properties of GLS assumes that such moment conditions apply for the relevant case at hand. For completeness, as the third GLS assumption, $X_i$ and $u_i$ are simply assumed to satisfy suitable moment conditions.

The fourth GLS assumption is that $X$ has full column rank; that is, the regressors are not perfectly multicollinear.

The GLS assumptions are summarized in Key Concept 18.4.

We consider GLS estimation in two cases. In the first case, $\mathbf{\Omega}(\mathbf{X})$ is known. In the second case, the functional form of $\mathbf{\Omega}(\mathbf{X})$ is known up to some parameters that can be estimated. To simplify notation, we refer to the function $\mathbf{\Omega}(\mathbf{X})$ as the matrix $\mathbf{\Omega}$, so the dependence of $\mathbf{\Omega}$ on $\mathbf{X}$ is implicit.

## GLS When $\Omega$ Is Known

When $\mathbf{\Omega}$ is known, the GLS estimator uses $\mathbf{\Omega}$ to transform the regression model to one with errors that satisfy the Gauss–Markov conditions. Specifically, let $\mathbf{F}$ be a matrix square root of $\mathbf{\Omega}^{-1}$; that is, let $\mathbf{F}$ be a matrix that satisfies $\mathbf{F}'\mathbf{F} = \mathbf{\Omega}^{-1}$ (see Appendix 18.1). A property of $\mathbf{F}$ is that $\mathbf{F\Omega F}' = \mathbf{I}_n$. Now premultiply both sides of Equation (18.4) by $\mathbf{F}$ to obtain

$$\widetilde{Y} = \widetilde{X}\boldsymbol{\beta} + \widetilde{U}, \tag{18.42}$$

where $\widetilde{\mathbf{Y}} = \mathbf{FY}$, $\widetilde{\mathbf{X}} = \mathbf{FX}$, and $\widetilde{\mathbf{U}} = \mathbf{FU}$.

The key insight of GLS is that, under the four GLS assumptions, the Gauss–Markov assumptions hold for the transformed regression in Equation (18.42). That is, by transforming all the variables by the inverse of the matrix square root of $\mathbf{\Omega}$, the regression errors in the transformed regression have a conditional mean of zero and a covariance matrix that equals the identity matrix. To show this mathematically, first note that $E(\widetilde{U}|\widetilde{X}) = E(\mathbf{FU}|\mathbf{FX}) = \mathbf{F}E(\mathbf{U}|\mathbf{FX}) = \mathbf{0}_n$ by the first GLS assumption [Equation (18.40)]. In addition, $E(\widetilde{U}\widetilde{U}'|\widetilde{X}) = E[(\mathbf{FU})(\mathbf{FU})'|\mathbf{FX}] = \mathbf{F}E(\mathbf{UU}'|\mathbf{FX})\mathbf{F}' = \mathbf{F\Omega F}' = \mathbf{I}_n$, where the second equality follows because $(\mathbf{FU})' = \mathbf{U}'\mathbf{F}'$ and the final equality follows from the definition of $\mathbf{F}$. It follows that the transformed regression model in Equation (18.42) satisfies the Gauss–Markov conditions in Key Concept 18.3.

The GLS estimator, $\widetilde{\boldsymbol{\beta}}^{GLS}$, is the OLS estimator of $\boldsymbol{\beta}$ in Equation (18.42); that is, $\widetilde{\boldsymbol{\beta}}^{GLS} = (\widetilde{X}'\widetilde{X})^{-1}(\widetilde{X}'\widetilde{Y})$. Because the transformed regression model satisfies the Gauss–Markov conditions, the GLS estimator is the best conditionally unbiased estimator that is linear in $\widetilde{Y}$. But because $\widetilde{Y} = \mathbf{FY}$ and $\mathbf{F}$ is (here) assumed to be known, and because $\mathbf{F}$ is invertible (because $\mathbf{\Omega}$ is positive definite), the class of estimators that are linear in $\widetilde{Y}$ is the same as the class of estimators that are linear in $\mathbf{Y}$. Thus the OLS estimator of $\boldsymbol{\beta}$ in Equation (18.42) is also the best conditionally unbiased estimator among estimators that are linear in $\mathbf{Y}$. In other words, under the GLS assumptions, the GLS estimator is BLUE.

The GLS estimator can be expressed directly in terms of $\boldsymbol{\Omega}$, so in principle there is no need to compute the square root matrix $\boldsymbol{F}$. Because $\widetilde{\boldsymbol{X}} = \boldsymbol{FX}$ and $\widetilde{\boldsymbol{Y}} = \boldsymbol{FY}$, $\widetilde{\boldsymbol{\beta}}^{GLS} = (\boldsymbol{X}'\boldsymbol{F}'\boldsymbol{FX})^{-1}(\boldsymbol{X}'\boldsymbol{F}'\boldsymbol{FY})$. But $\boldsymbol{F}'\boldsymbol{F} = \boldsymbol{\Omega}^{-1}$, so

$$\widetilde{\boldsymbol{\beta}}^{GLS} = (\boldsymbol{X}'\boldsymbol{\Omega}^{-1}\boldsymbol{X})^{-1}(\boldsymbol{X}'\boldsymbol{\Omega}^{-1}\boldsymbol{Y}). \tag{18.43}$$

In practice, $\boldsymbol{\Omega}$ is typically unknown, so the GLS estimator in Equation (18.43) typically cannot be computed and thus is sometimes called the **infeasible GLS** estimator. If, however, $\boldsymbol{\Omega}$ has a known functional form but the parameters of that function are unknown, then $\boldsymbol{\Omega}$ can be estimated and a feasible version of the GLS estimator can be computed.

## GLS When $\boldsymbol{\Omega}$ Contains Unknown Parameters

If $\boldsymbol{\Omega}$ is a known function of some parameters that in turn can be estimated, then these estimated parameters can be used to calculate an estimator of the covariance matrix $\boldsymbol{\Omega}$. For example, consider the time series application discussed following Equation (18.41), in which $\boldsymbol{\Omega}(\boldsymbol{X})$ does not depend on $\boldsymbol{X}$, $\boldsymbol{\Omega}_{ii} = \sigma_u^2$, $\boldsymbol{\Omega}_{ij} = \rho\sigma_u^2$ for $|i - j| = 1$, and $\boldsymbol{\Omega}_{ij} = 0$ for $|i - j| > 1$. Then $\boldsymbol{\Omega}$ has two unknown parameters, $\sigma_u^2$ and $\rho$. These parameters can be estimated using the residuals from a preliminary OLS regression; specifically, $\sigma_u^2$ can be estimated by $s_{\hat{u}}^2$ and $\rho$ can be estimated by the sample correlation between all neighboring pairs of OLS residuals. These estimated parameters can in turn be used to compute an estimator of $\boldsymbol{\Omega}$, $\hat{\boldsymbol{\Omega}}$.

In general, suppose that you have an estimator $\hat{\boldsymbol{\Omega}}$ of $\boldsymbol{\Omega}$. Then the GLS estimator based on $\hat{\boldsymbol{\Omega}}$ is

$$\hat{\boldsymbol{\beta}}^{GLS} = (\boldsymbol{X}'\hat{\boldsymbol{\Omega}}^{-1}\boldsymbol{X})^{-1}(\boldsymbol{X}'\hat{\boldsymbol{\Omega}}^{-1}\boldsymbol{Y}). \tag{18.44}$$

The GLS estimator in Equation (18.44) is sometimes called the **feasible GLS** estimator because it can be computed if the covariance matrix contains some unknown parameters that can be estimated.

## The Zero Conditional Mean Assumption and GLS

For the OLS estimator to be consistent, the first least squares assumption must hold; that is, $E(u_i|\boldsymbol{X}_i)$ must be zero. In contrast, the first GLS assumption is that $E(u_i|\boldsymbol{X}_1, \ldots, \boldsymbol{X}_n) = 0$. In other words, the first OLS assumption is that the error for the $i^{\text{th}}$ observation has a conditional mean of zero, given the values of the regressors for that observation, whereas the first GLS assumption is that $u_i$ has a conditional mean of zero, given the values of the regressors for *all* observations.

As discussed in Section 18.1, the assumptions that $E(u_i|X_i) = 0$ and that sampling is i.i.d. together imply that $E(u_i|X_1, \ldots, X_n) = 0$. Thus, when sampling is i.i.d. so that GLS is WLS, the first GLS assumption is implied by the first least squares assumption in Key Concept 18.1.

When sampling is not i.i.d., however, the first GLS assumption is not implied by the assumption that $E(u_i|X_i) = 0$; that is, the first GLS assumption is stronger. Although the distinction between these two conditions might seem slight, it can be very important in applications to time series data. This distinction is discussed in Section 15.5 in the context of whether the regressor is "past and present" exogenous or "strictly" exogenous; the assumption that $E(u_i|X_1, \ldots, X_n) = 0$ corresponds to strict exogeneity. Here, we discuss this distinction at a more general level using matrix notation. To do so, we focus on the case that $U$ is homoskedastic, $\Omega$ is known, and $\Omega$ has nonzero off-diagonal elements.

### *The role of the first GLS assumption.*

To see the source of the difference between these assumptions, it is useful to contrast the consistency arguments for GLS and OLS.

We first sketch the argument for the consistency of the GLS estimator in Equation (18.43). Substituting Equation (18.4) into Equation (18.43), we have $\widetilde{\boldsymbol{\beta}}^{GLS} = \boldsymbol{\beta} + (X'\Omega^{-1}X/n)^{-1}(X'\Omega^{-1}U/n)$. Under the first GLS assumption, $E(X'\Omega^{-1}U) = E[X'\Omega^{-1}E(U|X)] = \mathbf{0}_n$. If in addition the variance of $X'\Omega^{-1}U/n$ tends to zero and $X'\Omega^{-1}X/n \xrightarrow{p} \widetilde{Q}$, where $\widetilde{Q}$ is some invertible matrix, then $\widetilde{\boldsymbol{\beta}}^{GLS} \xrightarrow{p} \boldsymbol{\beta}$. Critically, when $\Omega$ has off-diagonal elements, the term $X'\Omega^{-1}U = \sum_{i=1}^{n}\sum_{j=1}^{n}X_i(\Omega^{-1})_{ij}u_j$ involves products of $X_i$ and $u_j$ for different $i, j$, where $(\Omega^{-1})_{ij}$ denotes the $(i, j)$ element of $\Omega^{-1}$. Thus, for $X'\Omega^{-1}U$ to have a mean of zero, it is not enough that $E(u_i|X_i) = 0$; rather $E(u_i|X_j)$ must equal zero for all $i, j$ pairs corresponding to nonzero values of $(\Omega^{-1})_{ij}$. Depending on the covariance structure of the errors, only some of or all the elements of $(\Omega^{-1})_{ij}$ might be nonzero. For example, if $u_i$ follows a first order autoregression (as discussed in Section 15.5), the only nonzero elements $(\Omega^{-1})_{ij}$ are those for which $|i - j| \leq 1$. In general, however, all the elements of $\Omega^{-1}$ can be nonzero, so in general for $X'\Omega^{-1}U/n \xrightarrow{p} \mathbf{0}_{(k+1)\times 1}$ (and thus for $\widetilde{\boldsymbol{\beta}}^{GLS}$ to be consistent) we need that $E(U|X) = \mathbf{0}_n$; that is, the first GLS assumption must hold.

In contrast, recall the argument that the OLS estimator is consistent. Rewrite Equation (18.14) as $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + (X'X/n)^{-1}\frac{1}{n}\sum_{i=1}^{n}X_iu_i$. If $E(u_i | X_i) = 0$, then the term $\frac{1}{n}\sum_{i=1}^{n}X_iu_i$ has mean zero, and if this term has a variance that tends to zero, it converges in probability to zero. If in addition $X'X/n \xrightarrow{p} Q_X$, then $\hat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$.

### *Is the first GLS assumption restrictive?*

The first GLS assumption requires that the errors for the $i^{\text{th}}$ observation be uncorrelated with the regressors for all other

observations. This assumption is dubious in some time series applications. This issue is discussed in Section 15.6 in the context of an empirical example, the relationship between the change in the price of a contract for future delivery of frozen orange concentrate and the weather in Florida. As explained there, the error term in the regression of price changes on the weather is plausibly uncorrelated with current and past values of the weather, so the first OLS assumption holds. However, this error term is plausibly correlated with future values of the weather, so the first GLS assumption does *not* hold.

This example illustrates a general phenomenon in economic time series data that arises when the value of a variable today is set in part based on expectations of the future: Those future expectations typically imply that the error term today depends on a forecast of the regressor tomorrow, which in turn is correlated with the actual value of the regressor tomorrow. For this reason, the first GLS assumption is in fact much stronger than the first OLS assumption. Accordingly, in some applications with economic time series data the GLS estimator is not consistent even though the OLS estimator is.

## 18.7   Instrumental Variables and Generalized Method of Moments Estimation

This section provides an introduction to the theory of instrumental variables (IV) estimation and the asymptotic distribution of IV estimators. It is assumed throughout that the IV regression assumptions in Key Concepts 12.3 and 12.4 hold and, moreover, that the instruments are strong. These assumptions apply to cross-sectional data with i.i.d. observations. Under certain conditions the results derived in this section are applicable to time series data as well, and the extension to time series data is briefly discussed at the end of this section. All asymptotic results in this section are developed under the assumption of strong instruments.

This section begins by presenting the IV regression model, the two stage least squares (TSLS) estimator, and its asymptotic distribution in the general case of heteroskedasticity, all in matrix form. It is next shown that, in the special case of homoskedasticity, the TSLS estimator is asymptotically efficient among the class of IV estimators in which the instruments are linear combinations of the exogenous variables. Moreover, the *J*-statistic has an asymptotic chi-squared distribution in which the degrees of freedom equal the number of overidentifying restrictions. This section concludes with a discussion of efficient IV estimation and the test of overidentifying restrictions when the errors are heteroskedastic—a situation in which the efficient IV estimator is known as the efficient generalized method of moments (GMM) estimator.

## The IV Estimator in Matrix Form

In this section, we let $X$ denote the $n \times (k + r + 1)$ matrix of the regressors in the equation of interest, so $X$ contains the included endogenous regressors (the $X$'s in Key Concept 12.1) *and* the included exogenous regressors (the $W$'s in Key Concept 12.1). That is, in the notation of Key Concept 12.1, the $i$th row of $X$ is $X_i' = (1 \quad X_{1i} \quad X_{2i} \quad \ldots \quad X_{ki} \quad W_{1i} \quad W_{2i} \quad \ldots \quad W_{ri})$. Also, let $Z$ denote the $n \times (m + r + 1)$ matrix of all the exogenous regressors, both those included in the equation of interest (the $W$'s) *and* those excluded from the equation of interest (the instruments). That is, in the notation of Key Concept 12.1, the $i$th row of $Z$ is $Z_i' = (1 \quad Z_{1i} \quad Z_{2i} \quad \ldots \quad Z_{mi} \quad W_{1i} \quad W_{2i} \quad \ldots \quad W_{ri})$.

With this notation, the IV regression model of Key Concept 12.1, written in matrix form, is

$$Y = X\boldsymbol{\beta} + U, \tag{18.45}$$

where $U$ is the $n \times 1$ vector of errors in the equation of interest, with $i$th element $u_i$.

The matrix $Z$ consists of all the exogenous regressors, so under the IV regression assumptions in Key Concept 12.4,

$$E(Z_i u_i) = 0 \quad \text{(instrument exogeneity).} \tag{18.46}$$

Because there are $k$ included endogenous regressors, the first stage regression consists of $k$ equations.

***The TSLS estimator.*** The TSLS estimator is the instrumental variables estimator in which the instruments are the predicted values of $X$ based on OLS estimation of the first stage regression. Let $\hat{X}$ denote this matrix of predicted values so that the $i$th row of $\hat{X}$ is $(\hat{X}_{1i} \quad \hat{X}_{2i} \quad \ldots \quad \hat{X}_{ki} \quad W_{1i} \quad W_{2i} \quad \ldots \quad W_{ri})$, where $\hat{X}_{1i}$ is the predicted value from the regression of $X_{1i}$ on $Z$, and so forth. Because the $W$'s are contained in $Z$, the predicted value from a regression of $W_{1i}$ on $Z$ is just $W_{1i}$, and so forth, so $\hat{X} = P_Z X$, where $P_Z = Z(Z'Z)^{-1}Z'$ [see Equation (18.27)]. Accordingly, the TSLS estimator is

$$\hat{\boldsymbol{\beta}}^{TSLS} = (\hat{X}'\hat{X})^{-1}\hat{X}'Y. \tag{18.47}$$

Because $\hat{X} = P_Z X$, $\hat{X}'\hat{X} = X'P_Z X$, and $\hat{X}'Y = X'P_Z Y$, the TSLS estimator can be rewritten as

$$\hat{\boldsymbol{\beta}}^{TSLS} = (X'P_Z X)^{-1}X'P_Z Y. \tag{18.48}$$

## Asymptotic Distribution of the TSLS Estimator

Substituting Equation (18.45) into Equation (18.48), rearranging, and multiplying by $\sqrt{n}$ yields the expression for the centered and scaled TSLS estimator:

$$\sqrt{n}(\hat{\boldsymbol{\beta}}^{TSLS} - \boldsymbol{\beta}) = \left(\frac{\boldsymbol{X'P_ZX}}{n}\right)^{-1} \frac{\boldsymbol{X'P_ZU}}{\sqrt{n}}$$

$$= \left[\frac{\boldsymbol{X'Z}}{n}\left(\frac{\boldsymbol{Z'Z}}{n}\right)^{-1}\frac{\boldsymbol{Z'X}}{n}\right]^{-1}\left[\frac{\boldsymbol{X'Z}}{n}\left(\frac{\boldsymbol{Z'Z}}{n}\right)^{-1}\frac{\boldsymbol{Z'U}}{\sqrt{n}}\right], \quad (18.49)$$

where the second equality uses the definition of $\boldsymbol{P_Z}$. Under the IV regression assumptions, $\boldsymbol{X'Z}/n \xrightarrow{p} \boldsymbol{Q_{XZ}}$ and $\boldsymbol{Z'Z}/n \xrightarrow{p} \boldsymbol{Q_{ZZ}}$, where $\boldsymbol{Q_{XZ}} = E(\boldsymbol{X}_i\boldsymbol{Z}_i')$ and $\boldsymbol{Q_{ZZ}} = E(\boldsymbol{Z}_i\boldsymbol{Z}_i')$. In addition, under the IV regression assumptions, $\boldsymbol{Z}_iu_i$ is i.i.d. with mean zero [Equation (18.46)] and a nonzero finite variance, so its sum, divided by $\sqrt{n}$, satisfies the conditions of the central limit theorem and

$$\boldsymbol{Z'U}/\sqrt{n} \xrightarrow{d} \boldsymbol{\Psi}_{ZU}, \text{ where } \boldsymbol{\Psi}_{ZU} \sim N(\boldsymbol{0}, \boldsymbol{H}), \boldsymbol{H} = E(\boldsymbol{Z}_i\boldsymbol{Z}_i'u_i^2) \quad (18.50)$$

and $\boldsymbol{\Psi}_{ZU}$ is $(m + r + 1) \times 1$.

Application of Equation (18.50) and of the limits $\boldsymbol{X'Z}/n \xrightarrow{p} \boldsymbol{Q_{XZ}}$ and $\boldsymbol{Z'Z}/n \xrightarrow{p} \boldsymbol{Q_{ZZ}}$ to Equation (18.49) yields the result that, under the IV regression assumptions, the TSLS estimator is asymptotically normally distributed:

$$\sqrt{n}\,(\hat{\boldsymbol{\beta}}^{TSLS} - \boldsymbol{\beta}) \xrightarrow{d} (\boldsymbol{Q_{XZ}Q_{ZZ}^{-1}Q_{ZX}})^{-1}\boldsymbol{Q_{XZ}Q_{ZZ}^{-1}\Psi}_{ZU} \sim N(\boldsymbol{0}, \boldsymbol{\Sigma}^{TSLS}), \quad (18.51)$$

where

$$\boldsymbol{\Sigma}^{TSLS} = (\boldsymbol{Q_{XZ}Q_{ZZ}^{-1}Q_{ZX}})^{-1}\boldsymbol{Q_{XZ}Q_{ZZ}^{-1}HQ_{ZZ}^{-1}Q_{ZX}}(\boldsymbol{Q_{XZ}Q_{ZZ}^{-1}Q_{ZX}})^{-1}, \quad (18.52)$$

where $\boldsymbol{H}$ is defined in Equation (18.50).

***Standard errors for TSLS.*** The formula in Equation (18.52) is daunting. Nevertheless, it provides a way to estimate $\boldsymbol{\Sigma}^{TSLS}$ by substituting sample moments for the population moments. The resulting variance estimator is

$$\hat{\boldsymbol{\Sigma}}^{TSLS} = (\hat{\boldsymbol{Q}}_{XZ}\hat{\boldsymbol{Q}}_{ZZ}^{-1}\hat{\boldsymbol{Q}}_{ZX})^{-1}\hat{\boldsymbol{Q}}_{XZ}\hat{\boldsymbol{Q}}_{ZZ}^{-1}\hat{\boldsymbol{H}}\hat{\boldsymbol{Q}}_{ZZ}^{-1}\hat{\boldsymbol{Q}}_{ZX}(\hat{\boldsymbol{Q}}_{XZ}\hat{\boldsymbol{Q}}_{ZZ}^{-1}\hat{\boldsymbol{Q}}_{ZX})^{-1}, \quad (18.53)$$

where $\hat{\boldsymbol{Q}}_{XZ} = \boldsymbol{X'Z}/n, \hat{\boldsymbol{Q}}_{ZZ} = \boldsymbol{Z'Z}/n, \hat{\boldsymbol{Q}}_{ZX} = \boldsymbol{Z'X}/n$, and

$$\hat{\boldsymbol{H}} = \frac{1}{n}\sum_{i=1}^{n}\boldsymbol{Z}_i\boldsymbol{Z}_i\,\hat{u}_i^2, \text{ where } \hat{\boldsymbol{U}} = \boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}}^{TSLS} \quad (18.54)$$

so that $\hat{U}$ is the vector of TSLS residuals and where $\hat{u}_i$ is the $i^{\text{th}}$ element of that vector (the TSLS residual for the $i^{\text{th}}$ observation).

The TSLS standard errors are the square roots of the diagonal elements of $\hat{\Sigma}^{TSLS}/n$.

## Properties of TSLS When the Errors Are Homoskedastic

If the errors are homoskedastic, then the TSLS estimator is asymptotically efficient among the class of IV estimators in which the instruments are linear combinations of the rows of $Z$. This result is the IV counterpart to the Gauss–Markov theorem and constitutes an important justification for using TSLS.

***The TSLS distribution under homoskedasticity.*** If the errors are homoskedastic, that is, if $E(u_i^2|Z_i) = \sigma_u^2$, then $H = E(Z_iZ_i'u_i^2) = E[E(Z_iZ_i'u_i^2|Z_i)] = E[Z_iZ_i'E(u_i^2|Z_i)] = Q_{ZZ}\sigma_u^2$. In this case, the variance of the asymptotic distribution of the TSLS estimator in Equation (18.52) simplifies to

$$\Sigma^{TSLS} = (Q_{XZ}Q_{ZZ}^{-1}Q_{ZX})^{-1}\sigma_u^2 \quad \text{(homoskedasticity only)}. \quad (18.55)$$

The homoskedasticity-only estimator of the TSLS variance matrix is

$$\widetilde{\Sigma}^{TSLS} = (\hat{Q}_{XZ}\hat{Q}_{ZZ}^{-1}\hat{Q}_{ZX})^{-1}\hat{\sigma}_u^2, \text{ where } \hat{\sigma}_u^2 = \frac{\hat{U}'\hat{U}}{n-k-r-1}$$
$$\text{(homoskedasticity only)}, \quad (18.56)$$

and the homoskedasticity-only TSLS standard errors are the square root of the diagonal elements of $\widetilde{\Sigma}^{TSLS}/n$.

***The class of IV estimators that use linear combinations of Z.*** The class of IV estimators that use linear combinations of $Z$ as instruments can be generated in two equivalent ways. Both start with the same moment equation: Under the assumption of instrument exogeneity, the errors $U = Y - X\beta$ are uncorrelated with the exogenous regressors; that is, at the true value of $\beta$, Equation (18.46) implies that

$$E[(Y - X\beta)'Z] = 0. \quad (18.57)$$

Equation (18.57) constitutes a system of $m + r + 1$ equations involving the $k + r + 1$ unknown elements of $\beta$. When $m > k$, these equations are redundant,

in the sense that all are satisfied at the true value of $\boldsymbol{\beta}$. When these population moments are replaced by their sample moments, the system of equations $(\boldsymbol{Y} - \boldsymbol{Xb})'\boldsymbol{Z} = 0$ can be solved for $\boldsymbol{b}$ when there is exact identification ($m = k$). This value of $\boldsymbol{b}$ is the IV estimator of $\boldsymbol{\beta}$. However, when there is overidentification ($m > k$), the system of equations typically cannot all be satisfied by the same value of $\boldsymbol{b}$ because of sampling variation—there are more equations than unknowns—and in general this system does not have a solution.

The first approach to the problem of estimating $\boldsymbol{\beta}$ when there is overidentification is to trade off the desire to satisfy each equation by minimizing a quadratic form involving all the equations. Specifically, let $\boldsymbol{A}$ be an $(m + r + 1) \times (m + r + 1)$ symmetric positive semidefinite weight matrix and let $\hat{\boldsymbol{\beta}}_A^{IV}$ denote the estimator that minimizes

$$\min_b (\boldsymbol{Y} - \boldsymbol{Xb})'\boldsymbol{ZAZ}'(\boldsymbol{Y} - \boldsymbol{Xb}). \tag{18.58}$$

The solution to this minimization problem is found by taking the derivative of the objective function with respect to $\boldsymbol{b}$, setting the result equal to zero, and rearranging. Doing so yields $\hat{\boldsymbol{\beta}}_A^{IV}$, the IV estimator based on the weight matrix $\boldsymbol{A}$:

$$\hat{\boldsymbol{\beta}}_A^{IV} = (\boldsymbol{X}'\boldsymbol{ZAZ}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{ZAZ}'\boldsymbol{Y}. \tag{18.59}$$

Comparison of Equations (18.59) and (18.48) shows that TSLS is the IV estimator with $\boldsymbol{A} = (\boldsymbol{Z}'\boldsymbol{Z})^{-1}$. That is, TSLS is the solution of the minimization problem in Equation (18.58) with $\boldsymbol{A} = (\boldsymbol{Z}'\boldsymbol{Z})^{-1}$.

The calculations leading to Equations (18.51) and (18.52), applied to $\hat{\boldsymbol{\beta}}_A^{IV}$, show that

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_A^{IV} - \boldsymbol{\beta}) \xrightarrow{d} N(\boldsymbol{0}, \boldsymbol{\Sigma}_A^{IV}), \text{ where}$$
$$\boldsymbol{\Sigma}_A^{IV} = (\boldsymbol{Q}_{XZ}\boldsymbol{A}\boldsymbol{Q}_{ZX})^{-1}\boldsymbol{Q}_{XZ}\boldsymbol{AHA}\boldsymbol{Q}_{ZX}(\boldsymbol{Q}_{XZ}\boldsymbol{A}\boldsymbol{Q}_{ZX})^{-1}. \tag{18.60}$$

The second way to generate the class of IV estimators that use linear combinations of $\boldsymbol{Z}$ is to consider IV estimators in which the instruments are $\boldsymbol{ZB}$, where $\boldsymbol{B}$ is an $(m + r + 1) \times (k + r + 1)$ matrix with full row rank. Then the system of $(k + r + 1)$ equations, $(\boldsymbol{Y} - \boldsymbol{Xb})'\boldsymbol{ZB} = 0$, can be solved uniquely for the $(k + r + 1)$ unknown elements of $\boldsymbol{b}$. Solving these equations for $\boldsymbol{b}$ yields $\hat{\boldsymbol{\beta}}^{IV} = (\boldsymbol{B}'\boldsymbol{Z}'\boldsymbol{X})^{-1}(\boldsymbol{B}'\boldsymbol{Z}'\boldsymbol{Y})$, and substitution of $\boldsymbol{B} = \boldsymbol{AZ}'\boldsymbol{X}$ into this expression yields Equation (18.59). Thus the two approaches to defining IV estimators that are linear combinations of the instruments yield the same family of IV estimators. It is conventional to work with the first approach, in which the IV estimator solves the quadratic minimization problem in Equation (18.58), and that is the approach taken here.

*Asymptotic efficiency of TSLS under homoskedasticity.* If the errors are homoskedastic, then $\boldsymbol{H} = \boldsymbol{Q_{ZZ}}\sigma_u^2$ and the expression for $\boldsymbol{\Sigma}_A^{IV}$ in Equation (18.60) becomes

$$\boldsymbol{\Sigma}_A^{IV} = (\boldsymbol{Q_{XZ}AQ_{ZX}})^{-1}\boldsymbol{Q_{XZ}AQ_{ZZ}AQ_{ZX}}(\boldsymbol{Q_{XZ}AQ_{ZX}})^{-1}\sigma_u^2. \tag{18.61}$$

To show that TSLS is asymptotically efficient among the class of estimators that are linear combinations of $\boldsymbol{Z}$ when the errors are homoskedastic, we need to show that, under homoskedasticity,

$$\boldsymbol{c}'\boldsymbol{\Sigma}_A^{IV}\boldsymbol{c} \geq \boldsymbol{c}'\boldsymbol{\Sigma}^{TSLS}\boldsymbol{c} \tag{18.62}$$

for all positive semidefinite matrices $\boldsymbol{A}$ and all $(k + r + 1) \times 1$ vectors $\boldsymbol{c}$, where $\boldsymbol{\Sigma}^{TSLS} = (\boldsymbol{Q_{XZ}Q_{ZZ}^{-1}Q_{ZX}})^{-1}\sigma_u^2$ [Equation (18.55)]. The inequality (18.62), which is proven in Appendix 18.6, is the same efficiency criterion as is used in the multivariate Gauss–Markov theorem in Key Concept 18.3. Consequently, TSLS is the efficient IV estimator under homoskedasticity, among the class of estimators in which the instruments are linear combinations of $\boldsymbol{Z}$.

*The J-statistic under homoskedasticity.* The *J*-statistic (Key Concept 12.6) tests the null hypothesis that all the overidentifying restrictions hold against the alternative that some or all of them do not hold.

The idea of the *J*-statistic is that, if the overidentifying restrictions hold, $u_i$ will be uncorrelated with the instruments and thus a regression of $\boldsymbol{U}$ on $\boldsymbol{Z}$ will have population regression coefficients that all equal zero. In practice, $\boldsymbol{U}$ is not observed, but it can be estimated by the TSLS residuals $\hat{\boldsymbol{U}}$, so a regression of $\hat{\boldsymbol{U}}$ on $\boldsymbol{Z}$ should yield statistically insignificant coefficients. Accordingly, the TSLS *J*-statistic is the homoskedasticity-only *F*-statistic testing the hypothesis that the coefficients on $\boldsymbol{Z}$ are all zero, in the regression of $\hat{\boldsymbol{U}}$ on $\boldsymbol{Z}$, multiplied by $(m + r + 1)$ so that the *F*-statistic is in its asymptotic chi-squared form.

An explicit formula for the *J*-statistic can be obtained using Equation (7.13) for the homoskedasticity-only *F*-statistic. The unrestricted regression is the regression of $\hat{\boldsymbol{U}}$ on the $m + r + 1$ regressors $\boldsymbol{Z}$, and the restricted regression has no regressors. Thus, in the notation of Equation (7.13), $SSR_{unrestricted} = \hat{\boldsymbol{U}}'\boldsymbol{M_Z}\hat{\boldsymbol{U}}$ and $SSR_{restricted} = \hat{\boldsymbol{U}}'\hat{\boldsymbol{U}}$, so $SSR_{restricted} - SSR_{unrestricted} = \hat{\boldsymbol{U}}'\hat{\boldsymbol{U}} - \hat{\boldsymbol{U}}'\boldsymbol{M_Z}\hat{\boldsymbol{U}} = \hat{\boldsymbol{U}}'\boldsymbol{P_Z}\hat{\boldsymbol{U}}$ and the *J*-statistic is

$$J = \frac{\hat{\boldsymbol{U}}'\boldsymbol{P_Z}\hat{\boldsymbol{U}}}{\hat{\boldsymbol{U}}'\boldsymbol{M_Z}\hat{\boldsymbol{U}}/(n - m - r - 1)}. \tag{18.63}$$

The method for computing the $J$-statistic described in Key Concept 12.6 entails testing only the hypothesis that the coefficients on the excluded instruments are zero. Although these two methods have different computational steps, they produce identical $J$-statistics (Exercise 18.14).

It is shown in Appendix 18.6 that, under the null hypothesis that $E(u_i \mathbf{Z}_i) = 0$,

$$J \xrightarrow{d} \chi^2_{m-k}. \tag{18.64}$$

## Generalized Method of Moments Estimation in Linear Models

If the errors are heteroskedastic, then the TSLS estimator is no longer efficient among the class of IV estimators that use linear combinations of $\mathbf{Z}$ as instruments. The efficient estimator in this case is known as the efficient generalized method of moments (GMM) estimator. In addition, if the errors are heteroskedastic, then the $J$-statistic as defined in Equation (18.63) no longer has a chi-squared distribution. However, an alternative formulation of the $J$-statistic, constructed using the efficient GMM estimator, does have a chi-squared distribution with $m - k$ degrees of freedom.

These results parallel the results for the estimation of the usual regression model with exogenous regressors and heteroskedastic errors: If the errors are heteroskedastic, then the OLS estimator is not efficient among estimators that are linear in $\mathbf{Y}$ (the Gauss–Markov conditions are not satisfied) and the homoskedasticity-only $F$-statistic no longer has an $F$ distribution, even in large samples. In the regression model with exogenous regressors and heteroskedasticity, the efficient estimator is weighted least squares; in the IV regression model with heteroskedasticity, the efficient estimator uses a different weighting matrix than TSLS, and the resulting estimator is the efficient GMM estimator.

*GMM estimation.* **Generalized method of moments (GMM)** estimation is a general method for the estimation of the parameters of linear or nonlinear models, in which the parameters are chosen to provide the best fit to multiple equations, each of which sets a sample moment to zero. These equations, which in the context of GMM are called moment conditions, typically cannot all be satisfied simultaneously. The GMM estimator trades off the desire to satisfy each of the equations by minimizing a quadratic objective function.

In the linear IV regression model with exogenous variables $\mathbf{Z}$, the class of GMM estimators consists of all the estimators that are solutions to the quadratic minimization problem in Equation (18.58). Thus the class of GMM estimators based on the full set of instruments $\mathbf{Z}$ with different-weight matrices $\mathbf{A}$ is the same as the class of IV estimators in which the instruments are linear combinations of $\mathbf{Z}$.

In the linear IV regression model, GMM is just another name for the class of estimators we have been studying—that is, estimators that solve Equation (18.58).

*The asymptotically efficient GMM estimator.* Among the class of GMM estimators, the **efficient GMM** estimator is the GMM estimator with the smallest asymptotic variance matrix [where the smallest variance matrix is defined as in Equation (18.62)]. Thus the result in Equation (18.62) can be restated as saying that TSLS is the efficient GMM estimator in the linear model when the errors are homoskedastic.

To motivate the expression for the efficient GMM estimator when the errors are heteroskedastic, recall that when the errors are homoskedastic, $H$ [the variance matrix of $Z_i u_i$; see Equation (18.50)] equals $Q_{ZZ} \sigma_u^2$, and the asymptotically efficient weight matrix is obtained by setting $A = (Z'Z)^{-1}$, which yields the TSLS estimator. In large samples, using the weight matrix $A = (Z'Z)^{-1}$ is equivalent to using $A = (Q_{ZZ} \sigma_u^2)^{-1} = H^{-1}$. This interpretation of the TSLS estimator suggests that, by analogy, the efficient IV estimator under heteroskedasticity can be obtained by setting $A = H^{-1}$ and solving

$$\min_b (Y - Xb)'ZH^{-1}Z'(Y - Xb). \tag{18.65}$$

This analogy is correct: The solution to the minimization problem in Equation (18.65) is the efficient GMM estimator. Let $\widetilde{\beta}^{Eff.GMM}$ denote the solution to the minimization problem in Equation (18.65). By Equation (18.59), this estimator is

$$\widetilde{\beta}^{Eff.GMM} = (X'ZH^{-1}Z'X)^{-1}X'ZH^{-1}Z'Y. \tag{18.66}$$

The asymptotic distribution of $\widetilde{\beta}^{Eff.GMM}$ is obtained by substituting $A = H^{-1}$ into Equation (18.60) and simplifying; thus

$$\sqrt{n}(\widetilde{\beta}^{Eff.GMM} - \beta) \xrightarrow{d} N(0, \Sigma^{Eff.GMM}),$$
$$\text{where } \Sigma^{Eff.GMM} = (Q_{XZ}H^{-1}Q_{ZX})^{-1}. \tag{18.67}$$

The result that $\widetilde{\beta}^{Eff.GMM}$ is the efficient GMM estimator is proven by showing that $c'\Sigma_A^{IV} c \geq c'\Sigma^{Eff.GMM} c$ for all vectors $c$, where $\Sigma_A^{IV}$ is given in Equation (18.60). The proof of this result is given in Appendix 18.6.

*Feasible efficient GMM estimation.* The GMM estimator defined in Equation (18.66) is not a feasible estimator because it depends on the unknown variance matrix $H$. However, a feasible efficient GMM estimator can be computed by substituting a consistent estimator of $H$ into the minimization problem of Equation (18.65) or, equivalently, by substituting a consistent estimator of $H$ into the formula for $\hat{\beta}^{Eff.GMM}$ in Equation (18.66).

The efficient GMM estimator can be computed in two steps. In the first step, estimate $\boldsymbol{\beta}$ using any consistent estimator. Use this estimator of $\boldsymbol{\beta}$ to compute the residuals from the equation of interest, and then use these residuals to compute an estimator of $\boldsymbol{H}$. In the second step, use this estimator of $\boldsymbol{H}$ to estimate the optimal weight matrix $\boldsymbol{H}^{-1}$ and to compute the efficient GMM estimator. To be concrete, in the linear IV regression model, it is natural to use the TSLS estimator in the first step and to use the TSLS residuals to estimate $\boldsymbol{H}$. If TSLS is used in the first step, then the feasible efficient GMM estimator computed in the second step is

$$\hat{\boldsymbol{\beta}}^{Eff.GMM} = (\boldsymbol{X}'\boldsymbol{Z}\hat{\boldsymbol{H}}^{-1}\boldsymbol{Z}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{Z}\hat{\boldsymbol{H}}^{-1}\boldsymbol{Z}'\boldsymbol{Y}, \qquad (18.68)$$

where $\hat{\boldsymbol{H}}$ is given in Equation (18.54).

Because $\hat{\boldsymbol{H}} \xrightarrow{p} \boldsymbol{H}, \sqrt{n}(\hat{\boldsymbol{\beta}}^{Eff.GMM} - \tilde{\boldsymbol{\beta}}^{Eff.GMM}) \xrightarrow{p} 0$ (Exercise 18.12), and

$$\sqrt{n}\,(\hat{\boldsymbol{\beta}}^{Eff.GMM} - \boldsymbol{\beta}) \xrightarrow{d} N(0, \boldsymbol{\Sigma}^{Eff.GMM}), \qquad (18.69)$$

where $\boldsymbol{\Sigma}^{Eff.GMM} = (\boldsymbol{Q}_{XZ}\boldsymbol{H}^{-1}\boldsymbol{Q}_{ZX})^{-1}$ [Equation (18.67)]. That is, the feasible two-step estimator $\hat{\boldsymbol{\beta}}^{Eff.GMM}$ in Equation (18.68) is, asymptotically, the efficient GMM estimator.

### *The heteroskedasticity-robust J-statistic.* The **heteroskedasticity-robust J-statistic**, also known as the **GMM J-statistic**, is the counterpart of the TSLS-based $J$-statistic, computed using the efficient GMM estimator and weight function. That is, the GMM $J$-statistic is given by

$$J^{GMM} = (\boldsymbol{Z}'\hat{\boldsymbol{U}}^{GMM})'\hat{\boldsymbol{H}}^{-1}(\boldsymbol{Z}'\hat{\boldsymbol{U}}^{GMM})/n, \qquad (18.70)$$

where $\hat{\boldsymbol{U}}^{GMM} = \boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}}^{Eff.GMM}$ are the residuals from the equation of interest, estimated by (feasible) efficient GMM, and $\hat{\boldsymbol{H}}^{-1}$ is the weight matrix used to compute $\hat{\boldsymbol{\beta}}^{Eff.GMM}$.

Under the null hypothesis $E(\boldsymbol{Z}_i u_i) = \boldsymbol{0}, J^{GMM} \xrightarrow{d} \chi^2_{m-k}$ (see Appendix 18.6).

### *GMM with time series data.* The results in this section were derived under the IV regression assumptions for cross–sectional data. In many applications, however, these results extend to time series applications of IV regression and GMM. Although a formal mathematical treatment of GMM with time series data is beyond the scope of this book (for such a treatment, see Hayashi, 2000, Chapter 6), we nevertheless will summarize the key ideas of GMM estimation with time series data. This summary assumes familiarity with the material in Chapters 14 and 15. For this discussion, it is assumed that the variables are stationary.

It is useful to distinguish between two types of applications: applications in which the error term $u_t$ is serially correlated and applications in which $u_t$ is serially uncorrelated. If the error term $u_t$ is serially correlated, then the asymptotic distribution of the GMM estimator continues to be normally distributed, but the formula for $H$ in Equation (18.50) is no longer correct. Instead, the correct expression for $H$ depends on the autocovariances of $Z_t u_t$ and is analogous to the formula given in Equation (15.14) for the variance of the OLS estimator when the error term is serially correlated. The efficient GMM estimator is still constructed using a consistent estimator of $H$; however, that consistent estimator must be computed using the HAC methods discussed in Chapter 15.

If the error term $u_t$ is not serially correlated, then HAC estimation of $H$ is unnecessary and the formulas presented in this section all extend to time series GMM applications. In modern applications to finance and macroeconometrics, it is common to encounter models in which the error term represents an unexpected or unforecastable disturbance, in which case the model implies that $u_t$ is serially uncorrelated. For example, consider a model with a single included endogenous variable and no included exogenous variables so that the equation of interest is $Y_t = \beta_0 + \beta_1 X_t + u_t$. Suppose that an economic theory implies that $u_t$ is unpredictable given past information. Then the theory implies the moment condition

$$E(u_t \mid Y_{t-1}, X_{t-1}, Z_{t-1}, Y_{t-2}, X_{t-2}, Z_{t-2}, \dots) = 0, \qquad (18.71)$$

where $Z_{t-1}$ is the lagged value of some other variable. The moment condition in Equation (18.71) implies that all the lagged variables $Y_{t-1}, X_{t-1}, Z_{t-1}, Y_{t-2}, X_{t-2}, Z_{t-2}, \dots$ are candidates for being valid instruments (they satisfy the exogeneity condition). Moreover, because $u_{t-1} = Y_{t-1} - \beta_0 - \beta_1 X_{t-1}$, the moment condition in Equation (18.71) is equivalent to $E(u_t \mid u_{t-1}, X_{t-1}, Z_{t-1}, u_{t-2}, X_{t-2}, Z_{t-2}, \dots) = 0$. Because $u_t$ is serially uncorrelated, HAC estimation of $H$ is unnecessary. The theory of GMM presented in this section, including efficient GMM estimation and the GMM $J$-statistic, therefore applies directly to time series applications with moment conditions of the form in Equation (18.71), under the hypothesis that the moment condition in Equation (18.71) is, in fact, correct.

## Summary

1. The linear multiple regression model in matrix form is $Y = X\beta + U$, where $Y$ is the $n \times 1$ vector of observations on the dependent variable, $X$ is the $n \times (k + 1)$ matrix of $n$ observations on the $k + 1$ regressors (including a constant), $\beta$ is the $k + 1$ vector of unknown parameters, and $U$ is the $n \times 1$ vector of error terms.

2. The OLS estimator is $\hat{\boldsymbol{\beta}} = (\boldsymbol{X'X})^{-1}\boldsymbol{X'Y}$. Under the first four least squares assumptions in Key Concept 18.1, $\hat{\boldsymbol{\beta}}$ is consistent and asymptotically normally distributed. If in addition the errors are homoskedastic, then the conditional variance of $\hat{\boldsymbol{\beta}}$ is $\text{var}(\hat{\boldsymbol{\beta}} \mid \boldsymbol{X}) = \sigma_u^2(\boldsymbol{X'X})^{-1}$.

3. General linear restrictions on $\boldsymbol{\beta}$ can be written as the $q$ equations $\boldsymbol{R\beta} = \boldsymbol{r}$, and this formulation can be used to test joint hypotheses involving multiple coefficients or to construct confidence sets for elements of $\boldsymbol{\beta}$.

4. When the regression errors are i.i.d. and normally distributed, conditional on $\boldsymbol{X}$, $\boldsymbol{\beta}$ has an exact normal distribution and the homoskedasticity-only $t$- and $F$-statistics have exact $t_{n-k-1}$ and $F_{q,\,n-k-1}$ distributions, respectively.

5. The Gauss–Markov theorem says that, if the errors are homoskedastic and conditionally uncorrelated across observations and if $E(u_i|\boldsymbol{X}) = 0$, the OLS estimator is efficient among linear conditionally unbiased estimators (that is, OLS is BLUE).

6. If the error covariance matrix $\boldsymbol{\Omega}$ is not proportional to the identity matrix, and if $\boldsymbol{\Omega}$ is known or can be estimated, then the GLS estimator is asymptotically more efficient than OLS. However, GLS requires that, in general, $u_i$ be uncorrelated with *all* observations on the regressors, not just with $\boldsymbol{X}_i$, as is required by OLS, an assumption that must be evaluated carefully in applications.

7. The TSLS estimator is a member of the class of GMM estimators of the linear model. In GMM, the coefficients are estimated by making the sample covariance between the regression error and the exogenous variables as small as possible—specifically, by solving $\min_b [(\boldsymbol{Y} - \boldsymbol{Xb})'\boldsymbol{Z}]\boldsymbol{A}[\boldsymbol{Z}'(\boldsymbol{Y} - \boldsymbol{Xb})]$, where $\boldsymbol{A}$ is a weight matrix. The asymptotically efficient GMM estimator sets $\boldsymbol{A} = [E(\boldsymbol{Z}_i\boldsymbol{Z}_i'u_i^2)]^{-1}$. When the errors are homoskedastic, the asymptotically efficient GMM estimator in the linear IV regression model is TSLS.

## Key Terms

Gauss–Markov conditions for multiple regression (720)

Gauss–Markov theorem for multiple regression (721)

generalized least squares (GLS) (723)

infeasible GLS (726)

feasible GLS (726)

generalized method of moments (GMM) (734)

efficient GMM (735)

heteroskedasticity-robust                          mean vector (750)

    *J*-statistic (736)                                    covariance matrix (750)

GMM *J*-statistic (736)

---

**MyEconLab Can Help You Get a Better Grade**

---

## Review the Concepts

**18.1** A researcher studying the relationship between earnings and gender for a group of workers specifies the regression model $Y_i = \beta_0 + X_{1i}\beta_1 + X_{2i}\beta_2 + u_i$, where $X_{1i}$ is a binary variable that equals 1 if the $i^{\text{th}}$ person is a female and $X_{2i}$ is a binary variable that equals 1 if the $i^{\text{th}}$ person is a male. Write the model in the matrix form of Equation (18.2) for a hypothetical set of $n = 5$ observations. Show that the columns of $X$ are linearly dependent so that $X$ does not have full rank. Explain how you would respecifiy the model to eliminate the perfect multicollinearity.

**18.2** You are analyzing a linear regression model with 500 observations and one regressor. Explain how you would construct a confidence interval for $\beta_1$ if:

  **a.** Assumptions #1 through #4 in Key Concept 18.1 are true, but you think Assumption #5 or #6 might not be true.

  **b.** Assumptions #1 through #5 are true, but you think Assumption #6 might not be true. (Give two ways to construct the confidence interval.)

  **c.** Assumptions #1 through #6 are true.

**18.3** Suppose that Assumptions #1 through #5 in Key Concept 18.1 are true but that Assumption #6 is not. Does the result in Equation (18.31) hold? Explain.

**18.4** Can you compute the BLUE estimator of $\boldsymbol{\beta}$ if Equation (18.41) holds and you do not know $\boldsymbol{\Omega}$? What if you know $\boldsymbol{\Omega}$?

**18.5** Construct an example of a regression model that satisfies the assumption $E(u_i \mid X_i) = 0$ but for which $E(U \mid X) \neq \mathbf{0}_n$.

## Exercises

**18.1** Consider the population regression of test scores against income and the square of income in Equation (8.1).

    **a.** Write the regression in Equation (8.1) in the matrix form of Equation (18.5). Define $Y$, $X$, $U$, and $\beta$.

    **b.** Explain how to test the null hypothesis that the relationship between test scores and income is linear against the alternative that it is quadratic. Write the null hypothesis in the form of Equation (18.20). What are $R$, $r$, and $q$?

**18.2** Suppose that a sample of $n = 20$ households has the sample means and sample covariances below for a dependent variable and two regressors:

| | Sample Means | Sample Covariances | | |
| --- | --- | --- | --- | --- |
| | | $Y$ | $X_1$ | $X_2$ |
| $Y$ | 6.39 | 0.26 | 0.22 | 0.32 |
| $X_1$ | 7.24 | | 0.80 | 0.28 |
| $X_2$ | 4.00 | | | 2.40 |

    **a.** Calculate the OLS estimates of $\beta_0$, $\beta_1$, and $\beta_2$. Calculate $s_{\hat{u}}^2$. Calculate the $R^2$ of the regression.

    **b.** Suppose that all six assumptions in Key Concept 18.1 hold. Test the hypothesis that $\beta_1 = 0$ at the 5% significance level.

**18.3** Let $W$ be an $m \times 1$ vector with covariance matrix $\Sigma_W$, where $\Sigma_W$ is finite and positive definite. Let $c$ be a nonrandom $m \times 1$ vector and let $Q = c'W$.

    **a.** Show that $\text{var}(Q) = c'\Sigma_W c$.

    **b.** Suppose that $c \neq 0_m$. Show that $0 < \text{var}(Q) < \infty$.

**18.4** Consider the regression model $Y_i = \beta_0 + \beta_1 X_i + u_i$ from Chapter 4 and assume that the least squares assumptions in Key Concept 4.3 hold.

    **a.** Write the model in the matrix form given in Equations (18.2) and (18.4).

    **b.** Show that Assumptions #1 through #4 in Key Concept 18.1 are satisfied.

    **c.** Use the general formula for $\hat{\beta}$ in Equation (18.11) to derive the expressions for $\hat{\beta}_0$ and $\hat{\beta}_1$ given in Key Concept 4.2.

    **d.** Show that the (1, 1) element of $\Sigma_{\hat{\beta}}$ in Equation (18.13) is equal to the expression for $\sigma_{\hat{\beta}_0}^2$ given in Key Concept 4.4.

**18.5** Let $P_X$ and $M_X$ be as defined in Equations (18.24) and (18.25).

    **a.** Prove that $P_X M_X = \mathbf{0}_{n \times n}$ and that $P_X$ and $M_X$ are idempotent.

    **b.** Derive Equations (18.27) and (18.28).

    **c.** Show that rank$(P_X) = k + 1$ and rank$(M_X) = n - k - 1$. [*Hint:* First solve Exercise 18.10 and then use the fact that trace$(AB) = $ trace$(BA)$ for conformable matrices $A$ and $B$.]

**18.6** Consider the regression model in matrix form, $Y = X\beta + W\gamma + U$, where $X$ is an $n \times k_1$ matrix of regressors and $W$ is an $n \times k_2$ matrix of regressors. Then, as shown in Exercise 18.17, the OLS estimator $\hat{\beta}$ can be expressed

$$\hat{\beta} = (X'M_W X)^{-1}(X'M_W Y).$$

Now let $\hat{\beta}_1^{BV}$ be the "binary variable" fixed effects estimator computed by estimating Equation (10.11) by OLS and let $\hat{\beta}_1^{DM}$ be the "de-meaning" fixed effects estimator computed by estimating Equation (10.14) by OLS, in which the entity-specific sample means have been subtracted from $X$ and $Y$. Use the expression for $\hat{\beta}$ given above to prove that $\hat{\beta}_1^{BV} = \hat{\beta}_1^{DM}$. [*Hint*: Write Equation (10.11) using a full set of fixed effects, $D1_i, D2_i, \ldots, Dn_i$ and no constant term. Include all of the fixed effects in $W$. Write out the matrix $M_W X$.]

**18.7** Consider the regression model $Y_i = \beta_1 X_i + \beta_2 W_i + u_i$, where for simplicity the intercept is omitted and all variables are assumed to have a mean of zero. Suppose that $X_i$ is distributed independently of $(W_i, u_i)$ but $W_i$ and $u_i$ might be correlated and let $\hat{\beta}_1$ and $\hat{\beta}_2$ be the OLS estimators for this model. Show that

    **a.** Whether or not $W_i$ and $u_i$ are correlated, $\hat{\beta}_1 \xrightarrow{p} \beta_1$.

    **b.** If $W_i$ and $u_i$ are correlated, then $\hat{\beta}_2$ is inconsistent.

    **c.** Let $\hat{\beta}_1^r$ be the OLS estimator from the regression of $Y$ on $X$ (the restricted regression that excludes $W$). Will $\hat{\beta}_1$ have a smaller asymptotic variance than $\hat{\beta}_1^r$, allowing for the possibility that $W_i$ and $u_i$ are correlated? Explain.

**18.8** Consider the regression model $Y_i = \beta_0 + \beta_1 X_i + u_i$, where $u_1 = \tilde{u}_1$ and $u_i = 0.5u_{i-1} + \tilde{u}_i$ for $i = 2, 3, \ldots, n$. Suppose that $\tilde{u}_i$ are i.i.d. with mean 0 and variance 1 and are distributed independently of $X_j$ for all $i$ and $j$.

    **a.** Derive an expression for $E(UU') = \Omega$.

**b.** Explain how to estimate the model by GLS without explicitly inverting the matrix $\boldsymbol{\Omega}$. (*Hint:* Transform the model so that the regression errors are $\tilde{u}_1, \tilde{u}_2, \ldots, \tilde{u}_n$.)

**18.9** This exercise shows that the OLS estimator of a subset of the regression coefficients is consistent under the conditional mean independence assumption stated in Appendix 7.2. Consider the multiple regression model in matrix form $\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{W}\boldsymbol{\gamma} + \boldsymbol{U}$, where $\boldsymbol{X}$ and $\boldsymbol{W}$ are, respectively, $n \times k_1$ and $n \times k_2$ matrices of regressors. Let $\boldsymbol{X}_i'$ and $\boldsymbol{W}_i'$ denote the $i^{\text{th}}$ rows of $\boldsymbol{X}$ and $\boldsymbol{W}$ [as in Equation (18.3)]. Assume that (i) $E(u_i|\boldsymbol{X}_i, \boldsymbol{W}_i) = \boldsymbol{W}_i'\boldsymbol{\delta}$, where $\boldsymbol{\delta}$ is a $k_2 \times 1$ vector of unknown parameters; (ii) $(\boldsymbol{X}_i, \boldsymbol{W}_i, Y_i)$ are i.i.d.; (iii) $(\boldsymbol{X}_i, \boldsymbol{W}_i, u_i)$ have four finite, nonzero moments; and (iv) there is no perfect multicollinearity. These are Assumptions #1 through #4 of Key Concept 18.1, with the conditional mean independence assumption (i) replacing the usual conditional mean zero assumption.

**a.** Use the expression for $\hat{\boldsymbol{\beta}}$ given in Exercise 18.6 to write $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = (n^{-1}\boldsymbol{X}'\boldsymbol{M_W}\boldsymbol{X})^{-1}(n^{-1}\boldsymbol{X}'\boldsymbol{M_W}\boldsymbol{U})$.

**b.** Show that $n^{-1}\boldsymbol{X}'\boldsymbol{M_W}\boldsymbol{X} \xrightarrow{p} \boldsymbol{\Sigma}_{XX} - \boldsymbol{\Sigma}_{XW}\boldsymbol{\Sigma}_{WW}^{-1}\boldsymbol{\Sigma}_{WX}$, where $\boldsymbol{\Sigma}_{XX} = E(\boldsymbol{X}_i\boldsymbol{X}_i')$, $\boldsymbol{\Sigma}_{XW} = E(\boldsymbol{X}_i\boldsymbol{W}_i')$, and so forth. [The matrix $\boldsymbol{A}_n \xrightarrow{p} \boldsymbol{A}$ if $\boldsymbol{A}_{n,ij} \xrightarrow{p} \boldsymbol{A}_{ij}$ for all $i, j$, where $\boldsymbol{A}_{n,ij}$ and $\boldsymbol{A}_{ij}$ are the $(i, j)$ elements of $\boldsymbol{A}_n$ and $\boldsymbol{A}$.]

**c.** Show that assumptions (i) and (ii) imply that $E(\boldsymbol{U}|\boldsymbol{X}, \boldsymbol{W}) = \boldsymbol{W}\boldsymbol{\delta}$.

**d.** Use (c) and the law of iterated expectations to show that $n^{-1}\boldsymbol{X}'\boldsymbol{M_W}\boldsymbol{U} \xrightarrow{p} \boldsymbol{0}_{k_1 \times 1}$.

**e.** Use (a) through (d) to conclude that, under conditions (i) through (iv), $\hat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$.

**18.10** Let $\boldsymbol{C}$ be a symmetric idempotent matrix.

**a.** Show that the eigenvalues of $\boldsymbol{C}$ are either 0 or 1. (*Hint:* Note that $\boldsymbol{C}\boldsymbol{q} = \gamma\boldsymbol{q}$ implies $0 = \boldsymbol{C}\boldsymbol{q} - \gamma\boldsymbol{q} = \boldsymbol{C}\boldsymbol{C}\boldsymbol{q} - \gamma\boldsymbol{q} = \gamma\boldsymbol{C}\boldsymbol{q} - \gamma\boldsymbol{q} = \gamma^2\boldsymbol{q} - \gamma\boldsymbol{q}$ and solve for $\gamma$.)

**b.** Show that $\text{trace}(\boldsymbol{C}) = \text{rank}(\boldsymbol{C})$.

**c.** Let $\boldsymbol{d}$ be an $n \times 1$ vector. Show that $\boldsymbol{d}'\boldsymbol{C}\boldsymbol{d} \geq 0$.

**18.11** Suppose that $\boldsymbol{C}$ is an $n \times n$ symmetric idempotent matrix with rank $r$ and let $\boldsymbol{V} \sim N(\boldsymbol{0}_n, \boldsymbol{I}_n)$.

  **a.** Show that $C = AA'$, where $A$ is $n \times r$ with $A'A = I_r$. (*Hint*: $C$ is positive semidefinite and can be written as $Q\Lambda Q'$, as explained in Appendix 18.1.)

  **b.** Show that $A'V \sim N(\mathbf{0}_r, I_r)$.

  **c.** Show that $V'CV \sim \chi_r^2$.

**18.12 a.** Show that $\widetilde{\beta}^{Eff.GMM}$ is the efficient GMM estimator—that is, that $\widetilde{\beta}^{Eff.GMM}$ in Equation (18.66) is the solution to Equation (18.65).

  **b.** Show that $\sqrt{n}(\hat{\beta}^{Eff.GMM} - \widetilde{\beta}^{Eff.GMM}) \xrightarrow{p} 0$.

  **c.** Show that $J^{GMM} \xrightarrow{d} \chi_{m-k}^2$.

**18.13** Consider the problem of minimizing the sum of squared residuals, subject to the constraint that $Rb = r$, where $R$ is $q \times (k + 1)$ with rank $q$. Let $\widetilde{\beta}$ be the value of $b$ that solves the constrained minimization problem.

  **a.** Show that the Lagrangian for the minimization problem is $L(b, \gamma) = (Y - Xb)'\,(Y - Xb) + \gamma'(Rb - r)$, where $\gamma$ is a $q \times 1$ vector of Lagrange multipliers.

  **b.** Show that $\widetilde{\beta} = \hat{\beta} - (X'X)^{-1}R'[R(X'X)^{-1}R']^{-1}(R\hat{\beta} - r)$.

  **c.** Show that $(Y - X\widetilde{\beta})'(Y - X\widetilde{\beta}) - (Y - X\hat{\beta})(Y - X\hat{\beta}) = (R\hat{\beta} - r)'[R(X'X)^{-1}R']^{-1}(R\hat{\beta} - r)$.

  **d.** Show that $\widetilde{F}$ in Equation (18.36) is equivalent to the homoskedasticity-only $F$-statistic in Equation (7.13).

**18.14** Consider the regression model $Y = X\beta + U$. Partition $X$ as $[X_1 \quad X_2]$ and $\beta$ as $[\beta_1' \quad \beta_2']'$, where $X_1$ has $k_1$ columns and $X_2$ has $k_2$ columns. Suppose that $X_2'Y = \mathbf{0}_{k_2 \times 1}$. Let $R = [I_{k_1} \quad \mathbf{0}_{k_1 \times k_2}]$.

  **a.** Show that $\hat{\beta}'(X'X)\hat{\beta} = (R\hat{\beta})'[R(X'X)^{-1}R']^{-1}(R\hat{\beta})$.

  **b.** Consider the regression described in Equation (12.17). Let $W = [\mathbf{1} \quad W_1 \quad W_2 \quad \ldots \quad W_r]$, where $\mathbf{1}$ is an $n \times 1$ vector of ones, $W_1$ is the $n \times 1$ vector with $i^{th}$ element $W_{1i}$, and so forth. Let $\hat{U}^{TSLS}$ denote the vector of two-stage least squares residuals.

    **i.** Show that $W'\hat{U}^{TSLS} = 0$.

    **ii.** Show that the method for computing the $J$-statistic described in Key Concept 12.6 (using a homoskedasticity-only $F$-statistic) and the formula in Equation (18.63) produce the same value for the $J$-statistic. [*Hint*: Use the results in (a), (b, i), and Exercise 18.13.]

**18.15** (Consistency of clustered standard errors.) Consider the panel data model $Y_{it} = \beta X_{it} + \alpha_i + u_{it}$, where all variables are scalars. Assume that Assumptions #1, #2, and #4 in Key Concept 10.3 hold and strengthen Assumption #3 so that $X_{it}$ and $u_{it}$ have eight nonzero finite moments. Let $M = I_T - T^{-1}\boldsymbol{\iota}\boldsymbol{\iota}'$, where $\boldsymbol{\iota}$ is a $T \times 1$ vector of ones. Also let $\boldsymbol{Y}_i = (Y_{i1} \quad Y_{i2} \quad \cdots \quad Y_{iT})'$, $\boldsymbol{X}_i = (X_{i1} \quad X_{i2} \quad \cdots \quad X_{iT})'$, $\boldsymbol{u}_i = (u_{i1} \quad u_{i2} \quad \cdots \quad u_{iT})'$, $\widetilde{\boldsymbol{Y}}_i = M\boldsymbol{Y}_i$, $\widetilde{\boldsymbol{X}}_i = M\boldsymbol{X}_i$, and $\widetilde{\boldsymbol{u}}_i = M\boldsymbol{u}_i$. For the asymptotic calculations in this problem, suppose that $T$ is fixed and $n \longrightarrow \infty$.

**a.** Show that the fixed effects estimator of $\beta$ from Section 10.3 can be written as $\hat{\beta} = (\sum_{i=1}^{n}\widetilde{\boldsymbol{X}}_i'\widetilde{\boldsymbol{X}}_i)^{-1}\sum_{i=1}^{n}\widetilde{\boldsymbol{X}}_i'\widetilde{\boldsymbol{Y}}_i$.

**b.** Show that $\hat{\beta} - \beta = (\sum_{i=1}^{n}\widetilde{\boldsymbol{X}}_i'\widetilde{\boldsymbol{X}}_i)^{-1}\sum_{i=1}^{n}\widetilde{\boldsymbol{X}}_i'\boldsymbol{u}_i$. (*Hint:* $M$ is idempotent.)

**c.** Let $Q_{\widetilde{X}} = T^{-1}E(\widetilde{\boldsymbol{X}}_i'\widetilde{\boldsymbol{X}}_i)$ and $\hat{Q}_{\widetilde{X}} = \frac{1}{nT}\sum_{i=1}^{n}\sum_{t=1}^{T}\widetilde{X}_{it}^2$. Show that $\hat{Q}_{\widetilde{X}} \xrightarrow{p} Q_{\widetilde{X}}$.

**d.** Let $\eta_i = \widetilde{\boldsymbol{X}}_i'\boldsymbol{u}_i/\sqrt{T}$ and $\sigma_\eta^2 = \text{var}(\eta_i)$. Show that $\sqrt{\frac{1}{n}}\sum_{i=1}^{n}\eta_i \xrightarrow{d} N(0, \sigma_\eta^2)$.

**e.** Use your answers to (b) through (d) to prove Equation (10.25); that is, show that $\sqrt{nT}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \sigma_\eta^2/Q_{\widetilde{X}}^2)$.

**f.** Let $\widetilde{\sigma}_{\eta,clustered}^2$ be the infeasible clustered variance estimator, computed using the true errors instead of the residuals so that $\widetilde{\sigma}_{\eta,clustered}^2 = \frac{1}{nT}\sum_{i=1}^{n}(\widetilde{\boldsymbol{X}}_i'\boldsymbol{u}_i)^2$. Show that $\widetilde{\sigma}_{\eta,clustered}^2 \xrightarrow{p} \sigma_\eta^2$.

**g.** Let $\hat{\widetilde{\boldsymbol{u}}}_i = \widetilde{\boldsymbol{Y}}_i - \hat{\beta}\widetilde{\boldsymbol{X}}_i$ and $\hat{\sigma}_{\eta,\,clustered}^2 = \frac{n}{n-1}\frac{1}{nT}\sum_{i=1}^{n}(\widetilde{\boldsymbol{X}}_i'\hat{\widetilde{\boldsymbol{u}}}_i)^2$ [this is Equation (10.27) in matrix form]. Show that $\hat{\sigma}_{\eta,\,clustered}^2 \xrightarrow{p} \sigma_\eta^2$. [*Hint:* Use an argument like that used in Equation (17.16) to show that $\hat{\sigma}_{\eta,\,clustered}^2 - \widetilde{\sigma}_{\eta,\,clustered}^2 \xrightarrow{p} 0$ and then use your answer to (f).]

**18.16** This exercise takes up the problem of missing data discussed in Section 9.2. Consider the regression model $Y_i = X_i\beta + u_i, i = 1, \ldots, n$, where all variables are scalars and the constant term/intercept is omitted for convenience.

**a.** Suppose that the least squares assumptions in Key Concept 4.3 are satisfied. Show that the least squares estimator of $\beta$ is unbiased and consistent.

**b.** Now suppose that some of the observations are missing. Let $I_i$ denote a binary random variable that indicates the nonmissing observations; that is, $I_i = 1$ if observation $i$ is not missing and $I_i = 0$ if observation $i$ is missing. Assume that $\{I_i, X_i, u_i\}$ are i.i.d.

i. Show that the OLS estimator can be written as

$$\hat{\beta} = \left( \sum_{i=1}^{n} I_i X_i X_i' \right)^{-1} \left( \sum_{i=1}^{n} I_i X_i Y_i \right) = \beta + \left( \sum_{i=1}^{n} I_i X_i X_i' \right)^{-1} \left( \sum_{i=1}^{n} I_i X_i u_i \right).$$

ii. Suppose that data are missing, "completely at random," in the sense that $\Pr(I_i = 1 | X_i, u_i) = p$, where $p$ is a constant. Show that $\hat{\beta}$ is unbiased and consistent.

iii. Suppose that the probability that the $i^{th}$ observation is missing depends of $X_i$, but not on $u_i$; that is, $\Pr(I_i = 1 | X_i, u_i) = p(X_i)$. Show that $\hat{\beta}$ is unbiased and consistent.

iv. Suppose that the probability that the $i^{th}$ observation is missing depends on both $X_i$ and $u_i$; that is, $\Pr(I_i = 1 | X_i, u_i) = p(X_i, u_i)$. Is $\hat{\beta}$ unbiased? Is $\hat{\beta}$ consistent? Explain.

c. Suppose that $\beta = 1$ and that $X_i$ and $u_i$ are mutually independent standard normal random variables [so that both $X_i$ and $u_i$ are distributed $N(0, 1)$]. Suppose that $I_i = 1$ when $Y_i \geq 0$, but $I_i = 0$ when $Y_i < 0$. Is $\hat{\beta}$ unbiased? Is $\hat{\beta}$ consistent? Explain.

**18.17** Consider the regression model in matrix form $Y = X\beta + W\gamma + U$, where $X$ and $W$ are matrices of regressors and $\beta$ and $\gamma$ are vectors of unknown regression coefficients. Let $\tilde{X} = M_W X$ and $\tilde{Y} = M_W Y$, where $M_W = I - W(W'W)^{-1}W$.

a. Show that the OLS estimators of $\beta$ and $\gamma$ can be written as

$$\begin{bmatrix} \hat{\beta} \\ \hat{\gamma} \end{bmatrix} = \begin{bmatrix} X'X & X'W \\ W'X & W'W \end{bmatrix}^{-1} \begin{bmatrix} X'Y \\ W'Y \end{bmatrix}$$

b. Show that

$$\begin{bmatrix} X'X & X'W \\ W'X & W'W \end{bmatrix}^{-1}$$
$$= \begin{bmatrix} (X'M_W X)^{-1} & -(X'M_W X)^{-1}X'W(W'W)^{-1} \\ -(W'W)^{-1}W'X(X'M_W X)^{-1} & (W'W)^{-1} + (W'W)^{-1}W'X(X'M_W X)^{-1}X'W(W'W)^{-1} \end{bmatrix}.$$

(*Hint:* Show that the product of the two matrices is equal to the identity matrix.)

   **c.** Show that $\hat{\boldsymbol{\beta}} = (X'M_W X)^{-1}X'M_W Y$.

   **d.** The Frisch–Waugh theorem (Appendix 6.2) says that $\hat{\boldsymbol{\beta}} = (\widetilde{X}'\widetilde{X})^{-1}\widetilde{X}'\widetilde{Y}$. Use the result in (c) to prove the Frisch–Waugh theorem.

---

# 18.1 Summary of Matrix Algebra

This appendix summarizes vectors, matrices, and the elements of matrix algebra used in Chapter 1. The purpose of this appendix is to review some concepts and definitions from a course in linear algebra, not to replace such a course.

## Definitions of Vectors and Matrices

A **vector** is a collection of $n$ numbers or elements, collected either in a column (a **column vector**) or in a row (a **row vector**). The $n$-dimensional column vector $\boldsymbol{b}$ and the $n$-dimensional row vector $\boldsymbol{c}$ are

$$\boldsymbol{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \text{ and } \boldsymbol{c} = \begin{bmatrix} c_1 & c_2 & \cdots & c_n \end{bmatrix},$$

where $b_1$ is the first element of $\boldsymbol{b}$ and in general $b_i$ is the $i^{\text{th}}$ element of $\boldsymbol{b}$.

Throughout, a boldface symbol denotes a vector or matrix.

A **matrix** is a collection, or an array, of numbers or elements in which the elements are laid out in columns and rows. The dimension of a matrix is $n \times m$, where $n$ is the number of rows and $m$ is the number of columns. The $n \times m$ matrix $\boldsymbol{A}$ is

$$\boldsymbol{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{bmatrix},$$

where $a_{ij}$ is the $(i, j)$ element of $\boldsymbol{A}$, that is, $a_{ij}$ is the element that appears in the $i^{\text{th}}$ row and $j^{\text{th}}$ column. An $n \times m$ matrix consists of $n$ row vectors or, alternatively, of $m$ column vectors.

To distinguish one-dimensional numbers from vectors and matrices, a one-dimensional number is called a **scalar**.

## Types of Matrices

*Square, symmetric, and diagonal matrices.*  A matrix is said to be **square** if the number of rows equals the number of columns. A square matrix is said to be **symmetric** if its $(i, j)$ element equals its $(j, i)$ element. A **diagonal** matrix is a square matrix in which all the off-diagonal elements equal zero; that is, if the square matrix $A$ is diagonal, then $a_{ij} = 0$ for $i \neq j$.

*Special matrices.*  An important matrix is the **identity matrix**, $I_n$, which is an $n \times n$ diagonal matrix with ones on the diagonal. The **null matrix**, $\mathbf{0}_{n \times m}$, is the $n \times m$ matrix with all elements equal to zero.

*The transpose.*  The **transpose** of a matrix switches the rows and the columns. That is, the transpose of a matrix turns the $n \times m$ matrix $A$ into the $m \times n$ matrix, which is denoted by $A'$, where the $(i, j)$ element of $A$ becomes the $(j, i)$ element of $A'$; said differently, the transpose of the matrix $A$ turns the rows of $A$ into the columns of $A'$. If $a_{ij}$ is the $(i, j)$ element of $A$, then $A'$ (the transpose of $A$) is

$$A' = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{n1} \\ a_{12} & a_{22} & \cdots & a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1m} & a_{2m} & \cdots & a_{nm} \end{bmatrix}.$$

The transpose of a vector is a special case of the transpose of a matrix. Thus the transpose of a vector turns a column vector into a row vector; that is, if $b$ is an $n \times 1$ column vector, then its transpose is the $1 \times n$ row vector

$$b' = \begin{bmatrix} b_1 & b_2 & \cdots & b_n \end{bmatrix}.$$

The transpose of a row vector is a column vector.

## Elements of Matrix Algebra: Addition and Multiplication

*Matrix addition.*  Two matrices $A$ and $B$ that have the same dimensions (for example, that are both $n \times m$) can be added together. The sum of two matrices is the sum of their elements; that is, if $C = A + B$, then $c_{ij} = a_{ij} + b_{ij}$. A special case of matrix addition is vector addition: If $a$ and $b$ are both $n \times 1$ column vectors, then their sum $c = a + b$ is the element-wise sum; that is, $c_i = a_i + b_i$.

*Vector and matrix multiplication.*  Let $a$ and $b$ be two $n \times 1$ column vectors. Then the product of the transpose of $a$ (which is itself a row vector) with $b$ is $a'b = \sum_{i=1}^{n} a_i b_i$. Applying this definition with $b = a$ yields $a'a = \sum_{i=1}^{n} a_i^2$.

Similarly, the matrices $A$ and $B$ can be multiplied together if they are conformable—that is, if the number of columns of $A$ equals the number of rows of $B$. Specifically, suppose

that $A$ has dimension $n \times m$ and $B$ has dimension $m \times r$. Then the product of $A$ and $B$ is an $n \times r$ matrix, $C$; that is, $C = AB$, where the $(i, j)$ element of $C$ is $c_{ij} = \sum_{k=1}^{m} a_{ik}b_{kj}$. Said differently, the $(i, j)$ element of $AB$ is the product of multiplying the row vector that is the $i^{\text{th}}$ row of $A$ with the column vector that is the $j^{\text{th}}$ column of $B$.

The product of a scalar $d$ with the matrix $A$ has the $(i, j)$ element $da_{ij}$; that is, each element of $A$ is multiplied by the scalar $d$.

***Some useful properties of matrix addition and multiplication.***   Let $A$ and $B$ be matrices. Then:

    **a.**  $A + B = B + A$;

    **b.**  $(A + B) + C = A + (B + C)$;

    **c.**  $(A + B)' = A' + B'$;

    **d.**  If $A$ is $n \times m$, then $AI_m = A$ and $I_nA = A$;

    **e.**  $A(BC) = (AB)C$;

    **f.**  $(A + B)C = AC + BC$; and

    **g.**  $(AB)' = B'A'$.

In general, matrix multiplication does not commute; that is, in general $AB \neq BA$, although there are some special cases in which matrix multiplication commutes; for example, if $A$ and $B$ are both $n \times n$ diagonal matrices, then $AB = BA$.

## Matrix Inverse, Matrix Square Roots, and Related Topics

***The matrix inverse.***   Let $A$ be a square matrix. Assuming that it exists, the **inverse** of the matrix $A$ is defined as the matrix for which $A^{-1}A = I_n$. If in fact the inverse matrix $A^{-1}$ exists, then $A$ is said to be **invertible** or **nonsingular**. If both $A$ and $B$ are invertible, then $(AB)^{-1} = B^{-1}A^{-1}$.

***Positive definite and positive semidefinite matrices.***   Let $V$ be an $n \times n$ square matrix. Then $V$ is **positive definite** if $c'Vc > 0$ for all nonzero $n \times 1$ vectors $c$. Similarly, $V$ is **positive semidefinite** if $c'Vc \geq 0$ for all nonzero $n \times 1$ vectors $c$. If $V$ is positive definite, then it is invertible.

***Linear independence.***   The $n \times 1$ vectors $a_1$ and $a_2$ are **linearly independent** if there do not exist nonzero scalars $c_1$ and $c_2$ such that $c_1a_1 + c_2a_2 = 0_{n \times 1}$. More generally, the set of $k$ vectors $a_1, a_2, \ldots, a_k$ are linearly independent if there do not exist nonzero scalars $c_1, c_2, \ldots, c_k$ such that $c_1a_1 + c_2a_2 + \cdots + c_ka_k = 0_{n \times 1}$.

*The rank of a matrix.*  The **rank** of the $n \times m$ matrix $A$ is the number of linearly independent columns of $A$. The rank of $A$ is denoted $\text{rank}(A)$. If the rank of $A$ equals the number of columns of $A$, then $A$ is said to have full column rank. If the $n \times m$ matrix $A$ has full column rank, then there does not exist a nonzero $m \times 1$ vector $c$ such that $Ac = \mathbf{0}_{n \times 1}$. If $A$ is $n \times n$ with $\text{rank}(A) = n$, then $A$ is nonsingular. If the $n \times m$ matrix $A$ has full column rank, then $A'A$ is nonsingular.

*The matrix square root.*  Let $V$ be an $n \times n$ square symmetric positive definite matrix. The matrix square root of $V$ is defined to be an $n \times n$ matrix $F$ such that $F'F = V$. The matrix square root of a positive definite matrix will always exist, but it is not unique. The matrix square root has the property that $FV^{-1}F' = I_n$. In addition, the matrix square root of a positive definite matrix is invertible, so $F'^{-1}VF^{-1} = I_n$.

*Eigenvalues and eigenvectors.*  Let $A$ be an $n \times n$ matrix. If the $n \times 1$ vector $q$ and the scalar $\lambda$ satisfy $Aq = \lambda q$, where $q'q = 1$, then $\lambda$ is an **eigenvalue** of $A$, and $q$ is the **eigenvector** of $A$ associated with that eigenvalue. An $n \times n$ matrix has $n$ eigenvalues, which need not take on distinct values, and $n$ eigenvectors.

   If $V$ is an $n \times n$ symmetric positive definite matrix, then all the eigenvalues of $V$ are positive real numbers, and all the eigenvectors of $V$ are real. Also, $V$ can be written in terms of its eigenvalues and eigenvectors as $V = Q\Lambda Q'$, where $\Lambda$ is a diagonal $n \times n$ matrix with diagonal elements that equal the eigenvalues of $V$, and $Q$ is an $n \times n$ matrix consisting of the eigenvectors of $V$, arranged so that the $i^{\text{th}}$ column of $Q$ is the eigenvector corresponding to the eigenvalue that is the $i^{\text{th}}$ diagonal element of $\Lambda$. The eigenvectors are orthonormal, so $Q'Q = I_n$.

*Idempotent matrices.*  A matrix $C$ is idempotent if $C$ is square and $CC = C$. If $C$ is an $n \times n$ idempotent matrix that is also symmetric, then $C$ is positive semidefinite and $C$ has $r$ eigenvalues that equal 1 and $n - r$ eigenvalues that equal 0, where $r = \text{rank}(C)$ (Exercise 18.10).

**APPENDIX**

## 18.2  Multivariate Distributions

This appendix collects various definitions and facts about distributions of vectors of random variables. We start by defining the mean and covariance matrix of the $n$-dimensional random variable $V$. Next we present the multivariate normal distribution. Finally, we summarize some facts about the distributions of linear and quadratic functions of jointly normally distributed random variables.

## The Mean Vector and Covariance Matrix

The first and second moments of an $m \times 1$ vector of random variables, $V = (V_1 \quad V_2 \quad \cdots \quad V_m)'$, are summarized by its mean vector and covariance matrix.

Because $V$ is a vector, the vector of its means—that is, its **mean vector**—is $E(V) = \boldsymbol{\mu}_V$. The $i^{\text{th}}$ element of the mean vector is the mean of the $i^{\text{th}}$ element of $V$.

The **covariance matrix** of $V$ is the matrix consisting of the variance $\text{var}(V_i)$, $i = 1, \ldots, m$, along the diagonal and the $(i, j)$ off-diagonal elements $\text{cov}(V_i, V_j)$. In matrix form, the covariance matrix $\boldsymbol{\Sigma}_V$ is

$$\boldsymbol{\Sigma}_V = E[(V - \boldsymbol{\mu}_V)(V - \boldsymbol{\mu}_V)'] = \begin{bmatrix} \text{var}(V_1) & \cdots & \text{cov}(V_1, V_m) \\ \vdots & \ddots & \vdots \\ \text{cov}(V_m, V_1) & \cdots & \text{var}(V_m) \end{bmatrix}. \quad (18.72)$$

## The Multivariate Normal Distribution

The $m \times 1$ vector random variable $V$ has a multivariate normal distribution with mean vector $\boldsymbol{\mu}_V$ and covariance matrix $\boldsymbol{\Sigma}_V$ if it has the joint probability density function

$$f(V) = \frac{1}{\sqrt{(2\pi)^m \det(\boldsymbol{\Sigma}_V)}} \exp\left[ -\frac{1}{2}(V - \boldsymbol{\mu}_V)' \boldsymbol{\Sigma}_V^{-1}(V - \boldsymbol{\mu}_V) \right], \quad (18.73)$$

where $\det(\boldsymbol{\Sigma}_V)$ is the determinant of the matrix $\boldsymbol{\Sigma}_V$. The multivariate normal distribution is denoted $N(\boldsymbol{\mu}_V, \boldsymbol{\Sigma}_V)$.

An important fact about the multivariate normal distribution is that if two jointly normally distributed random variables are uncorrelated (equivalently, have a block-diagonal covariance matrix), then they are independently distributed. That is, let $V_1$ and $V_2$ be jointly normally distributed random variables with respective dimensions $m_1 \times 1$ and $m_2 \times 1$. Then if $\text{cov}(V_1, V_2) = E[(V_1 - \boldsymbol{\mu}_{V_1})(V_2 - \boldsymbol{\mu}_{V_2})'] = \mathbf{0}_{m_1 \times m_2}$, $V_1$ and $V_2$ are independent.

If $\{V_i\}$ are i.i.d. $N(0, \sigma_v^2)$, then $\boldsymbol{\Sigma}_V = \sigma_v^2 I_m$, and the multivariate normal distribution simplifies to the product of $m$ univariate normal densities.

## Distributions of Linear Combinations and Quadratic Forms of Normal Random Variables

Linear combinations of multivariate normal random variables are themselves normally distributed, and certain quadratic forms of multivariate normal random variables have a chi-squared distribution. Let $V$ be an $m \times 1$ random variable distributed $N(\boldsymbol{\mu}_V, \boldsymbol{\Sigma}_V)$, let $A$

and $\boldsymbol{B}$ be nonrandom $a \times m$ and $b \times m$ matrices, and let $\boldsymbol{d}$ be a nonrandom $a \times 1$ vector. Then

$$\boldsymbol{d} + \boldsymbol{A}\boldsymbol{V} \text{ is distributed } N(\boldsymbol{d} + \boldsymbol{A}\boldsymbol{\mu}_V, \boldsymbol{A}\boldsymbol{\Sigma}_V\boldsymbol{A}'); \tag{18.74}$$

$$\text{cov}(\boldsymbol{A}\boldsymbol{V}, \boldsymbol{B}\boldsymbol{V}) = \boldsymbol{A}\boldsymbol{\Sigma}_V\boldsymbol{B}'; \tag{18.75}$$

$$\text{if } \boldsymbol{A}\boldsymbol{\Sigma}_V\boldsymbol{B}' = \boldsymbol{0}_{a \times b}, \text{ then } \boldsymbol{A}\boldsymbol{V} \text{ and } \boldsymbol{B}\boldsymbol{V} \text{ are independently distributed; and} \tag{18.76}$$

$$(\boldsymbol{V} - \boldsymbol{\mu}_V)'\boldsymbol{\Sigma}_V^{-1}(\boldsymbol{V} - \boldsymbol{\mu}_V) \text{ is distributed } \chi_m^2. \tag{18.77}$$

Let $\boldsymbol{U}$ be an $m$-dimensional multivariate standard normal random variable with distribution $N(\boldsymbol{0}, \boldsymbol{I}_m)$. If $\boldsymbol{C}$ is symmetric and idempotent, then

$$\boldsymbol{U}'\boldsymbol{C}\boldsymbol{U} \text{ has a } \chi_r^2 \text{ distribution, where } r = \text{rank}(\boldsymbol{C}). \tag{18.78}$$

Equation (18.78) is proven as Exercise 18.11.

APPENDIX

# 18.3 Derivation of the Asymptotic Distribution of $\hat{\beta}$

This appendix provides the derivation of the asymptotic normal distribution of $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$ given in Equation (18.12). An implication of this result is that $\hat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$.

First consider the "denominator" matrix $\boldsymbol{X}'\boldsymbol{X}/n = \frac{1}{n}\sum_{i=1}^{n}\boldsymbol{X}_i\boldsymbol{X}_i'$ in Equation (18.15). The $(j, l)$ element of this matrix is $\frac{1}{n}\sum_{i=1}^{n} X_{ji}X_{li}$. By the second assumption in Key Concept 18.1, $\boldsymbol{X}_i$ is i.i.d., so $X_{ji}X_{li}$ is i.i.d. By the third assumption in Key Concept 18.1, each element of $\boldsymbol{X}_i$ has four moments, so, by the Cauchy–Schwarz inequality (Appendix 17.2), $X_{ji}X_{li}$ has two moments. Because $X_{ji}X_{li}$ is i.i.d. with two moments, $\frac{1}{n}\sum_{i=1}^{n} X_{ji}X_{li}$ obeys the law of large numbers, so $\frac{1}{n}\sum_{i=1}^{n} X_{ji}X_{li} \xrightarrow{p} E(X_{ji}X_{li})$. This is true for all the elements of $\boldsymbol{X}'\boldsymbol{X}/n$, so $\boldsymbol{X}'\boldsymbol{X}/n \xrightarrow{p} E(\boldsymbol{X}_i\boldsymbol{X}_i') = \boldsymbol{Q}_X$.

Next consider the "numerator" matrix in Equation (18.15), $\boldsymbol{X}'\boldsymbol{U}/\sqrt{n} = \sqrt{\frac{1}{n}}\sum_{i=1}^{n}\boldsymbol{V}_i$, where $\boldsymbol{V}_i = \boldsymbol{X}_iu_i$. By the first assumption in Key Concept 18.1 and the law of iterated expectations, $E(\boldsymbol{V}_i) = E[\boldsymbol{X}_iE(u_i|\boldsymbol{X}_i)] = \boldsymbol{0}_{k+1}$. By the second least squares assumption, $\boldsymbol{V}_i$ is i.i.d. Let $\boldsymbol{c}$ be a finite $k + 1$ dimensional vector. By the Cauchy–Schwarz inequality, $E[(\boldsymbol{c}'\boldsymbol{V}_i)^2] = E[(\boldsymbol{c}'\boldsymbol{X}_iu_i)^2] = E[(\boldsymbol{c}'\boldsymbol{X}_i)^2(u_i)^2] \leq \sqrt{E[(\boldsymbol{c}'\boldsymbol{X}_i)^4]E(u_i^4)}$, which is finite by the third least squares assumption. This is true for every such vector $\boldsymbol{c}$, so $E(\boldsymbol{V}_i\boldsymbol{V}_i') = \boldsymbol{\Sigma}_V$ is finite and, we assume, positive definite. Thus the multivariate central limit theorem of Key Concept 18.2 applies to $\sqrt{\frac{1}{n}}\sum_{i=1}^{n}\boldsymbol{V}_i = \frac{1}{\sqrt{n}}\boldsymbol{X}'\boldsymbol{U}$; that is,

$$\frac{1}{\sqrt{n}}\boldsymbol{X}'\boldsymbol{U} \xrightarrow{d} N(\boldsymbol{0}_{k+1}, \boldsymbol{\Sigma}_V). \tag{18.79}$$

The result in Equation (18.12) follows from Equations (18.15) and (18.79), the consistency of $X'X/n$, the fourth least squares assumption (which ensures that $(X'X)^{-1}$ exists), and Slutsky's theorem.

# 18.4 Derivations of Exact Distributions of OLS Test Statistics with Normal Errors

This appendix presents the proofs of the distributions under the null hypothesis of the homoskedasticity-only $t$-statistic in Equation (18.35) and the homoskedasticity-only $F$-statistic in Equation (18.37), assuming that all six assumptions in Key Concept 18.1 hold.

## Proof of Equation (18.35)

If (i) $Z$ has a standard normal distribution, (ii) $W$ has a $\chi_m^2$ distribution, and (iii) $Z$ and $W$ are independently distributed, then the random variable $Z/\sqrt{W/m}$ has the $t$-distribution with $m$ degrees of freedom (Appendix 17.1). To put $\tilde{t}$ in this form, notice that $\hat{\Sigma}_{\hat{\beta}} = (s_{\hat{u}}^2/\sigma_u^2)\Sigma_{\hat{\beta}|X}$. Then rewrite Equation (18.34) as

$$\tilde{t} = \frac{(\hat{\beta}_j - \beta_{j,0})/\sqrt{(\Sigma_{\hat{\beta}|X})_{jj}}}{\sqrt{W/(n-k-1)}},\tag{18.80}$$

where $W = (n-k-1)(s_{\hat{u}}^2/\sigma_u^2)$, and let $Z = (\hat{\beta}_j - \beta_{j,0})/\sqrt{(\Sigma_{\hat{\beta}|X})_{jj}}$ and $m = n-k-1$. With these definitions, $\tilde{t} = Z/\sqrt{W/m}$. Thus, to prove the result in Equation (18.35), we must show (i) through (iii) for these definitions of $Z$, $W$, and $m$.

i. An implication of Equation (18.30) is that, under the null hypothesis, $Z = (\hat{\beta}_j - \beta_{j,0})/\sqrt{(\Sigma_{\hat{\beta}|X})_{jj}}$ has an exact standard normal distribution, which shows (i).

ii. From Equation (18.31), $W$ is distributed as $\chi_{n-k-1}^2$, which shows (ii).

iii. To show (iii), it must be shown that $\hat{\beta}_j$ and $s_{\hat{u}}^2$ are independently distributed.

From Equations (18.14) and (18.29), $\hat{\beta} - \beta = (X'X)^{-1}X'U$ and $s_{\hat{u}}^2 = (M_X U)'(M_X U)/(n-k-1)$. Thus $\hat{\beta} - \beta$ and $s_{\hat{u}}^2$ are independent if $(X'X)^{-1}X'U$ and $M_X U$ are independent. Both $(X'X)^{-1}X'U$ and $M_X U$ are linear combinations of $U$, which has an $N(\mathbf{0}_{n\times1}, \sigma_u^2 I_n)$ distribution, conditional on $X$. But because $M_X X(X'X)^{-1} = \mathbf{0}_{n\times(k+1)}$ [Equation (18.26)], it follows that $(X'X)^{-1}X'U$ and $M_X U$ are independently distributed [Equation (18.76)]. Consequently, under all six assumptions in Key Concept 18.1,

$$\hat{\beta} \text{ and } s_{\hat{u}}^2 \text{ are independently distributed},\tag{18.81}$$

which shows (iii) and thus proves Equation (18.35).

## Proof of Equation (18.37)

The $F_{n_1, n_2}$ distribution is the distribution of $(W_1/n_1)/(W_2/n_2)$, where (i) $W_1$ is distributed $\chi^2_{n_1}$; (ii) $W_2$ is distributed $\chi^2_{n_2}$; and (iii) $W_1$ and $W_2$ are independently distributed (Appendix 17.1). To express $\widetilde{F}$ in this form, let $W_1 = (R\hat{\beta} - r)'[R(X'X)^{-1}R'\sigma^2_u]^{-1}(R\hat{\beta} - r)$ and $W_2 = (n - k - 1)s^2_{\hat{u}}/\sigma^2_u$ Substitution of these definitions into Equation (18.36) shows that $\widetilde{F} = (W_1/q)/[W_2/(n - k - 1)]$. Thus, by the definition of the $F$ distribution, $\widetilde{F}$ has an $F_{q, n-k-1}$ distribution if (i) through (iii) hold with $n_1 = q$ and $n_2 = n - k - 1$.

  i.  Under the null hypothesis, $R\hat{\beta} - r = R(\hat{\beta} - \beta)$. Because $\hat{\beta}$ has the conditional normal distribution in Equation (18.30) and because $R$ is a nonrandom matrix, $R(\hat{\beta} - \beta)$ is distributed $N(\mathbf{0}_{q\times 1}, R(X'X)^{-1}R'\sigma^2_u)$, conditional on $X$. Thus, by Equation (18.77) in Appendix 18.2, $(R\hat{\beta} - r)'[R(X'X)R'\sigma^2_u]^{-1}(R\hat{\beta} - r)$ is distributed $\chi^2_q$, proving (i).

  ii. Requirement (ii) is shown in Equation (18.31).

  iii. It has already been shown that $\hat{\beta} - \beta$ and $s^2_{\hat{u}}$ are independently distributed [Equation (18.81)]. It follows that $R\hat{\beta} - r$ and $s^2_{\hat{u}}$ are independently distributed, which in turn implies that $W_1$ and $W_2$ are independently distributed, proving (iii) and completing the proof.

**APPENDIX**

## 18.5 Proof of the Gauss–Markov Theorem for Multiple Regression

This appendix proves the Gauss–Markov theorem (Key Concept 18.3) for the multiple regression model. Let $\widetilde{\beta}$ be a linear conditionally unbiased estimator of $\beta$ so that $\widetilde{\beta} = A'Y$ and $E(\widetilde{\beta}\,|\,X) = \beta$, where $A$ is an $n \times (k + 1)$ matrix that can depend on $X$ and nonrandom constants. We show that $\mathrm{var}(c'\hat{\beta}) \le \mathrm{var}(c'\widetilde{\beta})$ for all $k + 1$ dimensional vectors $c$, where the inequality holds with equality only if $\widetilde{\beta} = \hat{\beta}$.

Because $\widetilde{\beta}$ is linear, it can be written as $\widetilde{\beta} = A'Y = A'(X\beta + U) = (A'X)\beta + A'U$. By the first Gauss–Markov condition, $E(U|X) = \mathbf{0}_{n\times 1}$, so $E(\widetilde{\beta}|X) = (A'X)\beta$, but because $\widetilde{\beta}$ is conditionally unbiased, $E(\widetilde{\beta}\,|\,X) = \beta = (A'X)\beta$, which implies that $A'X = I_{k+1}$. Thus $\widetilde{\beta} = \beta + A'U$, so $\mathrm{var}(\widetilde{\beta}|X) = \mathrm{var}(A'U|X) = E(A'UU'A|X) = A'E(UU'|X)A = \sigma^2_u A'A$, where the third equality follows because $A$ can depend on $X$ but not $U$, and the final equality follows from the second Gauss–Markov condition. That is, if $\widetilde{\beta}$ is linear and unbiased, then under the Gauss–Markov conditions,

$$A'X = I_{k+1} \text{ and } \mathrm{var}(\widetilde{\beta}|X) = \sigma^2_u A'A. \qquad (18.82)$$

The results in Equation (18.82) also apply to $\hat{\beta}$ with $A = \hat{A} = X(X'X)^{-1}$, where $(X'X)^{-1}$ exists by the third Gauss–Markov condition.

Now let $A = \hat{A} + D$ so that $D$ is the difference between the matrices $A$ and $\hat{A}$. Note that $\hat{A}'A = (X'X)^{-1}X'A = (X'X)^{-1}$ [by Equation (18.82)] and $\hat{A}'\hat{A} = (X'X)^{-1}X'X(X'X)^{-1} = (X'X)^{-1}$, so $\hat{A}'D = \hat{A}'(A - \hat{A}) = \hat{A}'A - \hat{A}'\hat{A} = \mathbf{0}_{(k+1)\times(k+1)}$. Substituting $A = \hat{A} + D$ into the formula for the conditional variance in Equation (18.82) yields

$$
\begin{aligned}
\text{var}(\tilde{\boldsymbol{\beta}}|X) &= \sigma_u^2(\hat{A} + D)'(\hat{A} + D) \\
&= \sigma_u^2[\hat{A}'\hat{A} + \hat{A}'D + D'\hat{A} + D'D] \\
&= \sigma_u^2(X'X)^{-1} + \sigma_u^2 D'D, \quad\quad (18.83)
\end{aligned}
$$

where the final equality uses the facts $\hat{A}'\hat{A} = (X'X)^{-1}$ and $\hat{A}'D = \mathbf{0}_{(k+1)\times(k+1)}$.

Because $\text{var}(\hat{\boldsymbol{\beta}}|X) = \sigma_u^2(X'X)^{-1}$, Equations (18.82) and (18.83) imply that $\text{var}(\tilde{\boldsymbol{\beta}}|X) - \text{var}(\hat{\boldsymbol{\beta}}|X) = \sigma_u^2 D'D$. The difference between the variances of the two estimators of the linear combination $c'\boldsymbol{\beta}$ thus is

$$
\text{var}(c'\tilde{\boldsymbol{\beta}}|X) - \text{var}(c'\hat{\boldsymbol{\beta}}|X) = \sigma_u^2 c'D'Dc \geq 0. \quad\quad (18.84)
$$

The inequality in Equation (18.84) holds for all linear combinations $c'\boldsymbol{\beta}$, and the inequality holds with equality for all nonzero $c$ only if $D = \mathbf{0}_{n\times(k+1)}$—that is, if $A = \hat{A}$ or, equivalently, $\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}$. Thus $c'\hat{\boldsymbol{\beta}}$ has the smallest variance of all linear conditionally unbiased estimators of $c'\boldsymbol{\beta}$; that is, the OLS estimator is BLUE.

# 18.6 Proof of Selected Results for IV and GMM Estimation

## The Efficiency of TSLS Under Homoskedasticity [Proof of Equation (18.62)]

When the errors $u_i$ are homoskedastic, the difference between $\Sigma_A^{IV}$ [Equation (18.61)] and $\Sigma^{TSLS}$ [Equation (18.55)] is given by

$$
\begin{aligned}
\Sigma_A^{IV} - \Sigma^{TSLS} &= (Q_{XZ}AQ_{ZX})^{-1}Q_{XZ}AQ_{ZZ}AQ_{ZX}(Q_{XZ}AQ_{ZX})^{-1}\sigma_u^2 - (Q_{XZ}Q_{ZZ}^{-1}Q_{ZX})^{-1}\sigma_u^2 \\
&= (Q_{XZ}AQ_{ZX})^{-1}Q_{XZ}A[Q_{ZZ} - Q_{ZX}(Q_{XZ}Q_{ZZ}^{-1}Q_{ZX})^{-1}Q_{XZ}]AQ_{ZX}(Q_{XZ}AQ_{ZX})^{-1}\sigma_u^2, \quad (18.85)
\end{aligned}
$$

where the second term in brackets in the second equality follows from $(Q_{XZ}AQ_{ZX})^{-1}Q_{XZ}AQ_{ZX} = I_{(k+r+1)}$. Let $F$ be the matrix square root of $Q_{ZZ}$, so $Q_{ZZ} = F'F$ and $Q_{ZZ}^{-1} = F^{-1}F'^{-1}$. [The latter equality follows from noting that $(F'F)^{-1} = F^{-1}F'^{-1}$ and $F'^{-1} = F^{-1'}$.] Then the final expression in Equation (18.85) can be rewritten to yield

$$
\begin{aligned}
\Sigma_A^{IV} - \Sigma^{TSLS} &= (Q_{XZ}AQ_{ZX})^{-1}Q_{XZ}AF'[I - F^{-1'}Q_{ZX}(Q_{XZ}F^{-1}F^{-1'}Q_{ZX})^{-1}Q_{XZ}F^{-1}] \\
&\quad \times FAQ_{ZX}(Q_{XZ}AQ_{ZX})^{-1}\sigma_u^2, \quad\quad (18.86)
\end{aligned}
$$