

Takeshi  
Amemiya

ADVANCED ECONOMETRICS

# Advanced Econometrics

---

Takeshi Amemiya

Basil Blackwell

© Takeshi Amemiya 1985

First published in 1985

This edition published in 1986 by Basil Blackwell Ltd,  
108 Cowley Road, Oxford OX4 1JF, UK

All rights reserved. Except for the quotation of short passages for the purposes of criticism and review, no part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior permission of the publisher.

Except in the United States of America, this book is sold subject to the condition that it shall not, by way of trade or otherwise, be lent, re-sold, hired out, or otherwise circulated without the publisher's prior consent in any form of binding or cover other than that in which it is published and without a similar condition including this condition being imposed on the subsequent purchaser.

*British Library Cataloguing in Publication Data*

Amemiya, Takeshi  
Advanced econometrics.  
1. Econometrics  
I. Title  
330'.028 HB139

ISBN 0-631-13345-3

Printed in Great Britain by TJ Press Ltd, Padstow

# Preface

---

This book is intended both as a reference book for professional econometricians and as a graduate textbook. If it is used as a textbook, the material contained in the book can be taught in a year-long course, as I have done at Stanford for many years. The prerequisites for such a course should be one year of calculus, one quarter or semester of matrix analysis, one year of intermediate statistical inference (see list of textbooks in note 1 of Chapter 3), and, preferably, knowledge of introductory or intermediate econometrics (say, at the level of Johnston, 1972). This last requirement is not necessary, but I have found in the past that a majority of economics students who take a graduate course in advanced econometrics do have knowledge of introductory or intermediate econometrics.

The main features of the book are the following: a thorough treatment of classical least squares theory (Chapter 1) and generalized least squares theory (Chapter 6); a rigorous discussion of large sample theory (Chapters 3 and 4); a detailed analysis of qualitative response models (Chapter 9), censored or truncated regression models (Chapter 10), and Markov chain and duration models (Chapter 11); and a discussion of nonlinear simultaneous equations models (Chapter 8).

The book presents only the fundamentals of time series analysis (Chapter 5 and a part of Chapter 6) because there are several excellent textbooks on the subject (see the references cited at the beginning of Chapter 5). In contrast, the models I discuss in the last three chapters have been used extensively in recent econometric applications but have not received in any textbook as complete a treatment as I give them here. Some instructors may wish to supplement my book with a textbook in time series analysis.

My discussion of linear simultaneous equations models (Chapter 7) is also brief. Those who wish to study the subject in greater detail should consult the references given in Chapter 7. I chose to devote more space to the discussion of nonlinear simultaneous equations models, which are still at an early stage of development and consequently have received only scant coverage in most textbooks.

In many parts of the book, and in all of Chapters 3 and 4, I have used the theorem-proof format and have attempted to develop all the mathematical results rigorously. However, it has not been my aim to present theorems in full mathematical generality. Because I intended this as a textbook rather than as a monograph, I chose assumptions that are relatively easy to understand and that lead to simple proofs, even in those instances where they could be relaxed. This will enable readers to understand the basic structure of each theorem and to generalize it for themselves depending on their needs and abilities. Many simple applications of theorems are given either in the form of examples in the text or in the form of exercises at the end of each chapter to bring out the essential points of each theorem.

Although this is a textbook in econometrics methodology, I have included discussions of numerous empirical papers to illustrate the practical use of theoretical results. This is especially conspicuous in the last three chapters of the book.

Too many people have contributed to the making of this book through the many revisions it has undergone to mention all their names. I am especially grateful to Trevor Breusch, Hidehiko Ichimura, Tom MaCurdy, Jim Powell, and Gene Savin for giving me valuable comments on the entire manuscript. I am also indebted to Carl Christ, Art Goldberger, Cheng Hsiao, Roger Koenker, Tony Lancaster, Chuck Manski, and Hal White for their valuable comments on parts of the manuscript. I am grateful to Colin Cameron, Tom Downes, Harry Paarsch, Aaron Han, and Choon Moon for proofreading and to the first three for correcting my English. In addition, Tom Downes and Choon Moon helped me with the preparation of the index. Dzung Pham has typed most of the manuscript through several revisions; her unflinching patience and good nature despite many hours of overtime work are much appreciated. David Criswell, Cathy Shimizu, and Bach-Hong Tran have also helped with the typing. The financial support of the National Science Foundation for the research that produced many of the results presented in the book is gratefully acknowledged. Finally, I am indebted to the editors of the *Journal of Economic Literature* for permission to include in Chapter 9 parts of my article entitled “Qualitative Response Models: A Survey” (*Journal of Economic Literature* 19:1483–1536, 1981) and to North-Holland Publishing Company for permission to use in Chapter 10 the revised version of my article entitled “Tobit Models: A Survey” (*Journal of Econometrics* 24:3–61, 1984).

# Advanced Econometrics

# Contents

---

<b>1</b>	<b>Classical Least Squares Theory</b>	<b>1</b>
<b>2</b>	<b>Recent Developments in Regression Analysis</b>	<b>45</b>
<b>3</b>	<b>Large Sample Theory</b>	<b>81</b>
<b>4</b>	<b>Asymptotic Properties of Extremum Estimators</b>	<b>105</b>
<b>5</b>	<b>Time Series Analysis</b>	<b>159</b>
<b>6</b>	<b>Generalized Least Squares Theory</b>	<b>181</b>
<b>7</b>	<b>Linear Simultaneous Equations Models</b>	<b>228</b>
<b>8</b>	<b>Nonlinear Simultaneous Equations Models</b>	<b>245</b>
<b>9</b>	<b>Qualitative Response Models</b>	<b>267</b>
<b>10</b>	<b>Tobit Models</b>	<b>360</b>
<b>11</b>	<b>Markov Chain and Duration Models</b>	<b>412</b>
	<b>Appendix 1 Useful Theorems in Matrix Analysis</b>	<b>459</b>
	<b>Appendix 2 Distribution Theory</b>	<b>463</b>
	<b>Notes</b>	<b>465</b>
	<b>References</b>	<b>475</b>
	<b>Name Index</b>	<b>505</b>
	<b>Subject Index</b>	<b>511</b>

# 1 Classical Least Squares Theory

---

In this chapter we shall consider the basic results of statistical inference in the classical linear regression model — the model in which the regressors are independent of the error term and the error term is serially uncorrelated and has a constant variance. This model is the starting point of the study; the models to be examined in later chapters are modifications of this one.

## 1.1 Linear Regression Model

In this section let us look at the reasons for studying the linear regression model and the method of specifying it. We shall start by defining Model 1, to be considered throughout the chapter.

### 1.1.1 Introduction

Consider a sequence of  $K$  random variables  $(y_t, x_{2t}, x_{3t}, \dots, x_{Kt})$ ,  $t = 1, 2, \dots, T$ . Define a  $T$ -vector  $\mathbf{y} = (y_1, y_2, \dots, y_T)'$ , a  $(K-1)$ -vector  $\mathbf{x}_t^* = (x_{2t}, x_{3t}, \dots, x_{Kt})'$ , and a  $[(K-1) \times T]$ -vector  $\mathbf{x}^* = (\mathbf{x}_1^{*'}, \mathbf{x}_2^{*'}, \dots, \mathbf{x}_T^{*'})'$ . Suppose for the sake of exposition that the joint density of the variables is given by  $f(\mathbf{y}, \mathbf{x}^*, \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  is a vector of unknown parameters. We are concerned with inference about the parameter vector  $\boldsymbol{\theta}$  on the basis of the observed vectors  $\mathbf{y}$  and  $\mathbf{x}^*$ .

In econometrics we are often interested in the conditional distribution of one set of random variables given another set of random variables; for example, the conditional distribution of consumption given income and the conditional distribution of quantities demanded given prices. Suppose we want to know the conditional distribution of  $\mathbf{y}$  given  $\mathbf{x}^*$ . We can write the joint density as the product of the conditional density and the marginal density as in

$$f(\mathbf{y}, \mathbf{x}^*, \boldsymbol{\theta}) = f(\mathbf{y}|\mathbf{x}^*, \boldsymbol{\theta}_1)f(\mathbf{x}^*, \boldsymbol{\theta}_2). \quad (1.1.1)$$

Regression analysis can be defined as statistical inferences on  $\boldsymbol{\theta}_1$ . For this purpose we can ignore  $f(\mathbf{x}^*, \boldsymbol{\theta}_2)$ , provided there is no relationship between  $\boldsymbol{\theta}_1$



and  $\theta_2$ . The vector  $y$  is called the vector of *dependent* or *endogenous* variables, and the vector  $x^*$  is called the vector of *independent* or *exogenous* variables.

In regression analysis we usually want to estimate only the first and second moments of the conditional distribution, rather than the whole parameter vector  $\theta_1$ . (In certain special cases the first two moments characterize  $\theta_1$  completely.) Thus we can define regression analysis as statistical inference on the conditional mean  $E(y|x^*)$  and the conditional variance-covariance matrix  $V(y|x^*)$ . Generally, these moments are nonlinear functions of  $x^*$ . However, in the present chapter we shall consider the special case in which  $E(y_i|x^*)$  is equal to  $E(y_i|x_i^*)$  and is a linear function of  $x_i^*$ , and  $V(y|x^*)$  is a constant times an identity matrix. Such a model is called the *classical* (or *standard*) linear regression model or the *homoscedastic* (meaning constant variance) linear regression model. Because this is the model to be studied in Chapter 1, let us call it simply Model 1.

### 1.1.2 Model 1

By writing  $x_t = (1, x_t^{*'})'$ , we can define Model 1 as follows. Assume

$$y_t = x_t' \beta + u_t, \quad t = 1, 2, \dots, T, \quad (1.1.2)$$

where  $y_t$  is a scalar observable random variable,  $\beta$  is a  $K$ -vector of unknown parameters,  $x_t$  is a  $K$ -vector of known constants such that  $\sum_{t=1}^T x_t x_t'$  is nonsingular, and  $u_t$  is a scalar, unobservable, random variable (called the error term or the disturbance) such that  $E u_t = 0$ ,  $V u_t = \sigma^2$  (another unknown parameter) for all  $t$ , and  $E u_t u_s = 0$  for  $t \neq s$ .

Note that we have assumed  $x^*$  to be a vector of known constants. This is essentially equivalent to stating that we are concerned only with estimating the conditional distribution of  $y$  given  $x^*$ . The most important assumption of Model 1 is the linearity of  $E(y_t|x_t^*)$ ; we therefore shall devote the next subsection to a discussion of the implications of that assumption. We have also made the assumption of homoscedasticity ( $V u_t = \sigma^2$  for all  $t$ ) and the assumption of no serial correlation ( $E u_t u_s = 0$  for  $t \neq s$ ), not because we believe that they are satisfied in most applications, but because they make a convenient starting point. These assumptions will be removed in later chapters.

We shall sometimes impose additional assumptions on Model 1 to obtain certain specific results. Notably, we shall occasionally make the assumption of serial independence of  $\{u_t\}$  or the assumption that  $u_t$  is normally distributed. In general, independence is a stronger assumption than no correlation, al-

though under normality the two concepts are equivalent. The additional assumptions will be stated whenever they are introduced into Model 1.

### 1.1.3 Implications of Linearity

Suppose random variables  $y_i$  and  $\mathbf{x}_i^*$  have finite second moments and their variance-covariance matrix is denoted by

$$V \begin{bmatrix} y_i \\ \mathbf{x}_i^* \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \sigma'_{12} \\ \sigma_{12} & \Sigma_{22} \end{bmatrix}.$$

Then we can always write

$$y_i = \beta_0 + \mathbf{x}_i^{*'} \beta_1 + v_i, \tag{1.1.3}$$

where  $\beta_1 = \Sigma_{22}^{-1} \sigma_{12}$ ,  $\beta_0 = E y_i - \sigma'_{12} \Sigma_{22}^{-1} E \mathbf{x}_i^*$ ,  $E v_i = 0$ ,  $V v_i = \sigma_1^2 - \sigma'_{12} \Sigma_{22}^{-1} \sigma_{12}$ , and  $E \mathbf{x}_i^{*'} v_i = 0$ . It is important to realize that Model 1 implies certain assumptions that (1.1.3) does not: (1.1.3) does not generally imply linearity of  $E(y_i | \mathbf{x}_i^*)$  because  $E(v_i | \mathbf{x}_i^*)$  may not generally be zero.

We call  $\beta_0 + \mathbf{x}_i^{*'} \beta_1$  in (1.1.3) the *best linear predictor* of  $y_i$  given  $\mathbf{x}_i^*$  because  $\beta_0$  and  $\beta_1$  can be shown to be the values of  $b_0$  and  $\mathbf{b}_1$  that minimize  $E(y_i - b_0 - \mathbf{x}_i^{*'} \mathbf{b}_1)^2$ . In contrast, the conditional mean  $E(y_i | \mathbf{x}_i^*)$  is called the *best predictor* of  $y_i$  given  $\mathbf{x}_i^*$  because  $E[y_i - E(y_i | \mathbf{x}_i^*)]^2 \leq E[y_i - g(\mathbf{x}_i^*)]^2$  for any function  $g$ .

The reader might ask why we work with eq. (1.1.2) rather than with (1.1.3). The answer is that (1.1.3) is so general that it does not allow us to obtain interesting results. For example, whereas the natural estimators of  $\beta_0$  and  $\beta_1$  can be defined by replacing the moments of  $y_i$  and  $\mathbf{x}_i^*$  that characterize  $\beta_0$  and  $\beta_1$  with their corresponding sample moments (they actually coincide with the least squares estimator), the mean of the estimator cannot be evaluated without specifying more about the relationship between  $\mathbf{x}_i^*$  and  $v_i$ .

How restrictive is the linearity of  $E(y_i | \mathbf{x}_i^*)$ ? It holds if  $y_i$  and  $\mathbf{x}_i^*$  are jointly normal or if  $y_i$  and  $\mathbf{x}_i^*$  are both scalar dichotomous (Bernoulli) variables.<sup>1</sup> But the linearity may not hold for many interesting distributions. Nevertheless, the linear assumption is not as restrictive as it may appear at first glance because  $\mathbf{x}_i^*$  can be variables obtained by transforming the original independent variables in various ways. For example, if the conditional mean of  $y_i$ , the supply of good, is a quadratic function of the price,  $p_i$ , we can put  $\mathbf{x}_i^{*'} = (p_i, p_i^2)'$ , thereby making  $E(y_i | \mathbf{x}_i^*)$  linear.

## 1.1.4 Matrix Notation

To facilitate the subsequent analysis, we shall write (1.1.2) in matrix notation as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad (1.1.4)$$

where  $\mathbf{y} = (y_1, y_2, \dots, y_T)'$ ,  $\mathbf{u} = (u_1, u_2, \dots, u_T)'$ , and  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)'$ . In other words,  $\mathbf{X}$  is the  $T \times K$  matrix, the  $t$ th row of which is  $\mathbf{x}_t'$ . The elements of the matrix  $\mathbf{X}$  are described as

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1K} \\ x_{21} & x_{22} & \dots & x_{2K} \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & & \cdot \\ x_{T1} & x_{T2} & \dots & x_{TK} \end{bmatrix}.$$

If we want to focus on the columns of  $\mathbf{X}$ , we can write  $\mathbf{X} = [\mathbf{x}_{(1)}, \mathbf{x}_{(2)}, \dots, \mathbf{x}_{(K)}]$ , where each  $\mathbf{x}_{(i)}$  is a  $T$ -vector. If there is no danger of confusing  $\mathbf{x}_{(i)}$  with  $\mathbf{x}_t$ , we can drop the parentheses and write simply  $\mathbf{x}_i$ . In matrix notation the assumptions on  $\mathbf{X}$  and  $\mathbf{u}$  can be stated as follows:  $\mathbf{X}'\mathbf{X}$  is nonsingular, which is equivalent to stating  $\text{rank}(\mathbf{X}) = K$  if  $T \geq K$ ;  $E\mathbf{u} = \mathbf{0}$ ; and  $E\mathbf{u}\mathbf{u}' = \sigma^2\mathbf{I}_T$ , where  $\mathbf{I}_T$  is the  $T \times T$  identity matrix. (Whenever the size of an identity matrix can be inferred from the context, we write it simply as  $\mathbf{I}$ .)

In the remainder of this chapter we shall no longer use the partition  $\boldsymbol{\beta}' = (\beta_0, \boldsymbol{\beta}_1')$ ; instead, the elements of  $\boldsymbol{\beta}$  will be written as  $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_K)'$ . Similarly, we shall not necessarily assume that  $\mathbf{x}_{(1)}$  is the vector of ones, although in practice this is usually the case. Most of our results will be obtained simply on the assumption that  $\mathbf{X}$  is a matrix of constants, without specifying specific values.

## 1.2 Theory of Least Squares

In this section we shall define the least squares estimator of the parameter  $\boldsymbol{\beta}$  in Model 1 and shall show that it is the best linear unbiased estimator. We shall also discuss estimation of the error variance  $\sigma^2$ .

1.2.1 Definition of Least Squares Estimators of  $\boldsymbol{\beta}$  and  $\sigma^2$ 

The least squares (LS) estimator  $\hat{\boldsymbol{\beta}}$  of the regression parameter  $\boldsymbol{\beta}$  in Model 1 is defined to be the value of  $\boldsymbol{\beta}$  that minimizes the sum of squared residuals<sup>2</sup>

$$\begin{aligned} S(\boldsymbol{\beta}) &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}. \end{aligned} \quad (1.2.1)$$

Putting the derivatives of  $S(\boldsymbol{\beta})$  with respect to  $\boldsymbol{\beta}$  equal to 0, we have

$$\frac{\partial S}{\partial \boldsymbol{\beta}} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{0}, \quad (1.2.2)$$

where  $\partial S/\partial \beta$  denotes the  $K$ -vector the  $i$ th element of which is  $\partial S/\partial \beta_i$ ,  $\beta_i$  being the  $i$ th element of  $\boldsymbol{\beta}$ . Solving (1.2.2) for  $\boldsymbol{\beta}$  gives

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}. \quad (1.2.3)$$

Clearly,  $S(\boldsymbol{\beta})$  attains the global minimum at  $\hat{\boldsymbol{\beta}}$ .

Let us consider the special case  $K = 2$  and  $\mathbf{x}'_i = (1, x_{2i})$  and represent each of the  $T$ -observations  $(y_i, x_{2i})$  by a point on the plane. Then, geometrically, the least squares estimates are the intercept and the slope of a line drawn in such a way that the sum of squares of the deviations between the points and the line is minimized in the direction of the  $y$ -axis. Different estimates result if the sum of squares of deviations is minimized in any other direction.

Given the least squares estimator  $\hat{\boldsymbol{\beta}}$ , we define

$$\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \quad (1.2.4)$$

and call it the vector of the *least squares residuals*. Using  $\hat{\mathbf{u}}$ , we can estimate  $\sigma^2$  by

$$\hat{\sigma}^2 = T^{-1}\hat{\mathbf{u}}'\hat{\mathbf{u}}, \quad (1.2.5)$$

called the least squares estimator of  $\sigma^2$ , although the use of the term *least squares* here is not as compelling as in the estimation of the regression parameters.

Using (1.2.4), we can write

$$\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{u}} = \mathbf{P}\mathbf{y} + \mathbf{M}\mathbf{y}, \quad (1.2.6)$$

where  $\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  and  $\mathbf{M} = \mathbf{I} - \mathbf{P}$ . Because  $\hat{\mathbf{u}}$  is orthogonal to  $\mathbf{X}$  (that is,  $\hat{\mathbf{u}}'\mathbf{X} = \mathbf{0}$ ), least squares estimation can be regarded as decomposing  $\mathbf{y}$  into two orthogonal components: a component that can be written as a linear combination of the column vectors of  $\mathbf{X}$  and a component that is orthogonal to  $\mathbf{X}$ . Alternatively, we can call  $\mathbf{P}\mathbf{y}$  the projection of  $\mathbf{y}$  onto the space spanned by the column vectors of  $\mathbf{X}$  and  $\mathbf{M}\mathbf{y}$  the projection of  $\mathbf{y}$  onto the space orthogonal to  $\mathbf{X}$ . Theorem 14 of Appendix 1 gives the properties of a projection matrix such as  $\mathbf{P}$  or  $\mathbf{M}$ . In the special case where both  $\mathbf{y}$  and  $\mathbf{X}$  are two-dimensional vectors

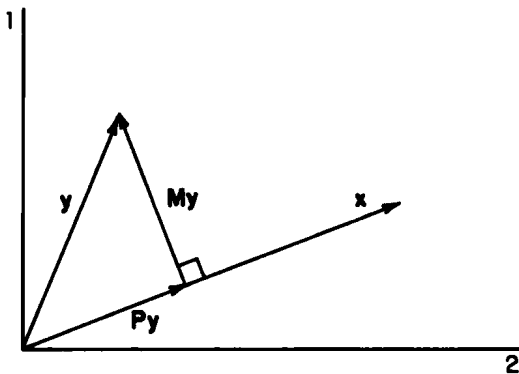


Figure 1.1 Orthogonal decomposition of  $y$

(that is,  $K = 1$  and  $T = 2$ ), the decomposition (1.2.6) can be illustrated as in Figure 1.1, where the vertical and horizontal axes represent the first and second observations, respectively, and the arrows represent vectors.

From (1.2.6) we obtain

$$y'y = y'Py + y'My. \tag{1.2.7}$$

The goodness of fit of the regression of  $y$  on  $X$  can be measured by the ratio  $y'Py/y'y$ , sometimes called  $R^2$ . However, it is more common to define  $R^2$  as the square of the sample correlation between  $y$  and  $Py$ :

$$R^2 = \frac{(y'LPy)^2}{y'Ly \cdot y'PLPy}, \tag{1.2.8}$$

where  $L = I_T - T^{-1}11'$  and  $1$  denotes the  $T$ -vector of ones. If we assume one of the columns of  $X$  is  $1$  (which is usually the case), we have  $LP = PL$ . Then we can rewrite (1.2.8) as

$$R^2 = \frac{y'LPLy}{y'Ly} = 1 - \frac{y'My}{y'Ly}. \tag{1.2.9}$$

Thus  $R^2$  can be interpreted as a measure of the goodness of fit of the regression of the deviations of  $y$  from its mean on the deviations of the columns of  $X$  from their means. (Section 2.1.4 gives a modification of  $R^2$  suggested by Theil, 1961.)

### 1.2.2 Least Squares Estimator of a Subset of $\beta$

It is sometimes useful to have an explicit formula for a subset of the least squares estimates  $\hat{\beta}$ . Suppose we partition  $\hat{\beta}' = (\hat{\beta}'_1, \hat{\beta}'_2)$ , where  $\hat{\beta}'_1$  is a  $K_1$ -vec-

tor and  $\hat{\beta}_2$  is a  $K_2$ -vector such that  $K_1 + K_2 = K$ . Partition  $X$  conformably as  $X = (X_1, X_2)$ . Then we can write  $X'X\hat{\beta} = X'y$  as

$$X_1'X_1\hat{\beta}_1 + X_1'X_2\hat{\beta}_2 = X_1'y \quad (1.2.10)$$

and

$$X_2'X_1\hat{\beta}_1 + X_2'X_2\hat{\beta}_2 = X_2'y. \quad (1.2.11)$$

Solving (1.2.11) for  $\hat{\beta}_2$  and inserting it into (1.2.10), we obtain

$$\hat{\beta}_1 = (X_1'M_2X_1)^{-1}X_1'M_2y, \quad (1.2.12)$$

where  $M_2 = I - X_2(X_2'X_2)^{-1}X_2'$ . Similarly,

$$\hat{\beta}_2 = (X_2'M_1X_2)^{-1}X_2'M_1y, \quad (1.2.13)$$

where  $M_1 = I - X_1(X_1'X_1)^{-1}X_1'$ .

In Model 1 we assume that  $X$  is of full rank, an assumption that implies that the matrices to be inverted in (1.2.12) and (1.2.13) are both nonsingular. Suppose for a moment that  $X_1$  is of full rank but that  $X_2$  is not. In this case  $\beta_2$  cannot be estimated, but  $\beta_1$  still can be estimated by modifying (1.2.12) as

$$\hat{\beta}_1 = (X_1'M_2^*X_1)^{-1}X_1'M_2^*y, \quad (1.2.14)$$

where  $M_2^* = I - X_2^*(X_2^{*'}X_2^*)^{-1}X_2^{*}$ , where the columns of  $X_2^*$  consist of a maximal number of linearly independent columns of  $X_2$ , provided that  $X_1'M_2^*X_1$  is nonsingular. (For the more general problem of estimating a linear combination of the elements of  $\beta$ , see Section 2.2.3.)

### 1.2.3 The Mean and Variance of $\hat{\beta}$ and $\hat{\sigma}^2$

Inserting (1.1.4) into (1.2.3), we have

$$\begin{aligned} \hat{\beta} &= (X'X)^{-1}X'y \\ &= \beta + (X'X)^{-1}X'u. \end{aligned} \quad (1.2.15)$$

Clearly,  $E\hat{\beta} = \beta$  by the assumptions of Model 1. Using the second line of (1.2.15), we can derive the variance-covariance matrix of  $\hat{\beta}$ :

$$\begin{aligned} V\hat{\beta} &= E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' \\ &= E(X'X)^{-1}X'uu'X(X'X)^{-1} \\ &= \sigma^2(X'X)^{-1}. \end{aligned} \quad (1.2.16)$$

From (1.2.3) and (1.2.4), we have  $\hat{u} = Mu$ , where  $M = I - X(X'X)^{-1}X'$ .

Using the properties of the projection matrix given in Theorem 14 of Appendix 1, we obtain

$$\begin{aligned}
 E\hat{\sigma}^2 &= T^{-1}E\mathbf{u}'\mathbf{M}\mathbf{u} && (1.2.17) \\
 &= T^{-1}E \operatorname{tr} \mathbf{M}\mathbf{u}\mathbf{u}' && \text{by Theorem 6 of Appendix 1} \\
 &= T^{-1}\sigma^2 \operatorname{tr} \mathbf{M} \\
 &= T^{-1}(T - K)\sigma^2 && \text{by Theorems 7 and 14 of Appendix 1,}
 \end{aligned}$$

which shows that  $\hat{\sigma}^2$  is a biased estimator of  $\sigma^2$ . We define the unbiased estimator of  $\sigma^2$  by

$$\tilde{\sigma}^2 = (T - K)^{-1}\hat{\mathbf{u}}'\hat{\mathbf{u}}. \quad (1.2.18)$$

We shall obtain the variance of  $\hat{\sigma}^2$  later, in Section 1.3, under the additional assumption that  $\mathbf{u}$  is normal.

The quantity  $V\hat{\beta}$  can be estimated by substituting either  $\hat{\sigma}^2$  or  $\tilde{\sigma}^2$  (defined above) for the  $\sigma^2$  that appears in the right-hand side of (1.2.16).

#### 1.2.4 Definition of Best

Before we prove that the least squares estimator is best linear unbiased, we must define the term *best*. First we shall define it for scalar estimators, then for vector estimators.

**DEFINITION 1.2.1.** Let  $\hat{\theta}$  and  $\theta^*$  be scalar estimators of a scalar parameter  $\theta$ . The estimator  $\hat{\theta}$  is said to be *at least as good as* (or *at least as efficient as*) the estimator  $\theta^*$  if  $E(\hat{\theta} - \theta)^2 \leq E(\theta^* - \theta)^2$  for all parameter values. The estimator  $\hat{\theta}$  is said to be *better* (or *more efficient*) than the estimator  $\theta^*$  if  $\hat{\theta}$  is at least as good as  $\theta^*$  and  $E(\hat{\theta} - \theta)^2 < E(\theta^* - \theta)^2$  for at least one parameter value. An estimator is said to be *best* (or *efficient*) in a class if it is better than any other estimator in the class.

The mean squared error is a reasonable criterion in many situations and is mathematically convenient. So, following the convention of the statistical literature, we have defined “better” to mean “having a smaller mean squared error.” However, there may be situations in which a researcher wishes to use other criteria, such as the mean absolute error.

**DEFINITION 1.2.2.** Let  $\hat{\theta}$  and  $\theta^*$  be estimators of a vector parameter  $\theta$ . Let  $\mathbf{A}$  and  $\mathbf{B}$  be their respective mean squared error matrices; that is,  $\mathbf{A} = E(\hat{\theta} - \theta)(\hat{\theta} - \theta)'$  and  $\mathbf{B} = E(\theta^* - \theta)(\theta^* - \theta)'$ . Then we say  $\hat{\theta}$  is *better* (or

more efficient) than  $\theta^*$  if

$$\mathbf{c}'(\mathbf{B} - \mathbf{A})\mathbf{c} \geq 0 \quad \text{for every vector } \mathbf{c} \quad \text{and every parameter value} \quad (1.2.19)$$

and

$$\mathbf{c}'(\mathbf{B} - \mathbf{A})\mathbf{c} > 0 \quad \text{for at least one value of } \mathbf{c} \quad \text{and} \quad (1.2.20) \\ \text{at least one value of the parameter.}$$

This definition of *better* clearly coincides with Definition 1.2.1 if  $\theta$  is a scalar.

In view of Definition 1.2.1, equivalent forms of statements (1.2.19) and (1.2.20) are statements (1.2.21) and (1.2.22):

$$\mathbf{c}'\hat{\theta} \quad \text{is at least as good as} \quad \mathbf{c}'\theta^* \quad \text{for every vector } \mathbf{c} \quad (1.2.21)$$

and

$$\mathbf{c}'\hat{\theta} \quad \text{is better than} \quad \mathbf{c}'\theta^* \quad \text{for at least one value of } \mathbf{c}. \quad (1.2.22)$$

Using Theorem 4 of Appendix 1, they also can be written as

$$\mathbf{B} \geq \mathbf{A} \quad \text{for every parameter value} \quad (1.2.23)$$

and

$$\mathbf{B} \neq \mathbf{A} \quad \text{for at least one parameter value.} \quad (1.2.24)$$

(Note that  $\mathbf{B} \geq \mathbf{A}$  means  $\mathbf{B} - \mathbf{A}$  is nonnegative definite and  $\mathbf{B} > \mathbf{A}$  means  $\mathbf{B} - \mathbf{A}$  is positive definite.)

We shall now prove the equivalence of (1.2.20) and (1.2.24). Because the phrase "for at least one parameter value" is common to both statements, we shall ignore it in the following proof. First, suppose (1.2.24) is not true. Then  $\mathbf{B} = \mathbf{A}$ . Therefore  $\mathbf{c}'(\mathbf{B} - \mathbf{A})\mathbf{c} = 0$  for every  $\mathbf{c}$ , a condition that implies that (1.2.20) is not true. Second, suppose (1.2.20) is not true. Then  $\mathbf{c}'(\mathbf{B} - \mathbf{A})\mathbf{c} = 0$  for every  $\mathbf{c}$  and every diagonal element of  $\mathbf{B} - \mathbf{A}$  must be 0 (choose  $\mathbf{c}$  to be the zero vector, except for 1 in the  $i$ th position). Also, the  $i, j$ th element of  $\mathbf{B} - \mathbf{A}$  is 0 (choose  $\mathbf{c}$  to be the zero vector, except for 1 in the  $i$ th and  $j$ th positions, and note that  $\mathbf{B} - \mathbf{A}$  is symmetric). Thus (1.2.24) is not true. This completes the proof.

Note that replacing  $\mathbf{B} \neq \mathbf{A}$  in (1.2.24) with  $\mathbf{B} > \mathbf{A}$ —or making the corresponding change in (1.2.20) or (1.2.22)—is unwise because we could not then rank the estimator with the mean squared error matrix

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



higher than the estimator with the mean squared error matrix

$$\begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}.$$

A problem with Definition 1.2.2 (more precisely, a problem inherent in the comparison of vector estimates rather than in this definition) is that often it does not allow us to say one estimator is either better or worse than the other. For example, consider

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{B} = \begin{bmatrix} 2 & 0 \\ 0 & \frac{1}{2} \end{bmatrix}. \quad (1.2.25)$$

Clearly, neither  $\mathbf{A} \geq \mathbf{B}$  nor  $\mathbf{B} \geq \mathbf{A}$ . In such a case one might compare the trace and conclude that  $\hat{\theta}$  is better than  $\theta^*$  because  $\text{tr } \mathbf{A} < \text{tr } \mathbf{B}$ . Another example is

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad \text{and} \quad \mathbf{B} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}. \quad (1.2.26)$$

Again, neither  $\mathbf{A} \geq \mathbf{B}$  nor  $\mathbf{B} \geq \mathbf{A}$ . If one were using the determinant as the criterion, one would prefer  $\hat{\theta}$  over  $\theta^*$  because  $\det \mathbf{A} < \det \mathbf{B}$ .

Note that  $\mathbf{B} \geq \mathbf{A}$  implies both  $\text{tr } \mathbf{B} \geq \text{tr } \mathbf{A}$  and  $\det \mathbf{B} \geq \det \mathbf{A}$ . The first follows from Theorem 7 and the second from Theorem 11 of Appendix 1. As these two examples show, neither  $\text{tr } \mathbf{B} \geq \text{tr } \mathbf{A}$  nor  $\det \mathbf{B} \geq \det \mathbf{A}$  implies  $\mathbf{B} \geq \mathbf{A}$ .

Use of the trace as a criterion has an obvious intuitive appeal, inasmuch as it is the sum of the individual variances. Justification for the use of the determinant involves more complicated reasoning. Suppose  $\hat{\theta} \sim N(\theta, \mathbf{V})$ , where  $\mathbf{V}$  is the variance-covariance matrix of  $\hat{\theta}$ . Then, by Theorem 1 of Appendix 2,  $(\hat{\theta} - \theta)' \mathbf{V}^{-1} (\hat{\theta} - \theta) \sim \chi_K^2$ , the chi-square distribution with  $K$  degrees of freedom,  $K$  being the number of elements of  $\theta$ . Therefore the  $(1 - \alpha)\%$  confidence ellipsoid for  $\theta$  is defined by

$$\{\theta | (\hat{\theta} - \theta)' \mathbf{V}^{-1} (\hat{\theta} - \theta) < \chi_K^2(\alpha)\}, \quad (1.2.27)$$

where  $\chi_K^2(\alpha)$  is the number such that  $P[\chi_K^2 \geq \chi_K^2(\alpha)] = \alpha$ . Then the volume of the ellipsoid (1.2.27) is proportional to the determinant of  $\mathbf{V}$ , as shown by Anderson (1958, p. 170).

A more intuitive justification for the determinant criterion is possible for the case in which  $\theta$  is a two-dimensional vector. Let the mean squared error matrix of an estimator  $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2)'$  be