

ECONOM \sim ETRICS
GREGORY C. CHOW

ECONOMETRICS

Gregory C. Chow

*Professor of Economics
Princeton University*

McGraw-Hill Book Company

New York St. Louis San Francisco Auckland Bogotá Hamburg
Johannesburg London Madrid Mexico Montreal New Delhi
Panama Paris São Paulo Singapore Sydney Tokyo Toronto

This book was set in Times Roman by Santype-Byrd.
The editors were Peter J. Dougherty and Scott Amerman;
the production supervisor was Leroy A. Young.
Halliday Lithograph Corporation was printer and binder.

ECONOMETRICS

Copyright © 1983 by McGraw-Hill, Inc. All rights reserved.
Printed in the United States of America. Except as permitted under the
United States Copyright Act of 1976, no part of this publication may be
reproduced or distributed in any form or by any means, or stored in a data
base or retrieval system, without the prior written permission of the
publisher.

234567890 HALHAL 89876543

ISBN 0-07-010847-1

Library of Congress Cataloging in Publication Data

Chow, Gregory C., date
Econometrics.

(Economics Handbook Series)

Includes bibliographies and index.

1. Econometrics. I. Title.

HB139.C483 1983 330'.028 82-20903

ISBN 0-07-010847-1

PREFACE

Since the 1970s, the development of econometrics has been so rapid that it has become difficult to find a textbook that treats the new topics and provides an updated perspective of the field. This text is intended to fill this need.

Most existing texts cover only the topics of Chaps. 1 to 5, the topics of Chaps. 6 to 12 being covered only infrequently. In order to incorporate new developments in a text, as it has been done in field after field for several decades, one must emphasize the basic ideas that underlie the subject. I believe from my teaching experience in Princeton that it is possible to cover the topics of seven chapters in about the same time and with the same effort on the part of students as previously spent in covering the topics of the first five chapters alone. This requires omitting nonessential details in, and improving the exposition of, previous material, as well as exposing the threads that connect all the material, old and new. After the new material has been integrated, the reader will gain a more up-to-date perspective of the subject.

This book is intended mainly for graduate students, although as the curriculum is constantly being upgraded, more and more advanced undergraduates will find the less theoretical parts readable. It assumes familiarity with the basic techniques of statistical inference and some exposure to matrix algebra. Naturally, more preparation in these two areas will make this text easier to read and simpler to comprehend. Chapter 1 can be used to test the reader's preparation. A reader with the minimum required preparation should find Secs. 1.1 to 1.5 easily comprehensible but may find the later starred sections difficult. In this case, I would recommend studying the nonstarred sections of Chap. 2 before returning to the more theoretical starred sections of Chap. 1. In fact, *an applications-*

oriented reader may skip the starred sections or read only the results therein without going through the proofs.

As another means of accommodating students with different backgrounds in statistics and mathematics, I have chosen to weave the required background material into the body of the text rather than presenting appendixes, which are usually dry and formal. This approach has the advantages of motivating the material and helping the student learn it while using it. Material that is used more often will be more firmly absorbed, as it should be.

The level of abstraction at which to present econometrics is a matter of taste. This book has interwoven empirical material to show students how empirical investigations in econometrics are conducted and to motivate the theoretical material. Some readers would prefer more empirical material than is actually presented and some less. Some would prefer having less theorem proving and others would prefer having more and at a higher level of mathematical abstraction. The level selected reflects the author's opinion concerning what the first-year graduate student in economics should be trained to do in the field of econometrics.

Although this text has a broader coverage than existing texts, it cannot possibly cover all important topics, as a textbook is not an encyclopedia. The selection of topics is also a matter of taste. Given the selection of topics, the choice of the tools and concepts to be presented within each topic has been guided by two criteria: the usefulness of the methods in applied work and the usefulness of the theoretical ideas in the further development of econometric methods. For example, the t test is useful in applied work, and the idea of maximum likelihood is useful for developing new estimators. Once the basic tools have been mastered, the student can apply them to study and solve new problems not covered in this text.

How this book should be used depends on the preparation of the students. As mentioned above, the starred sections can be omitted by applications-oriented readers. For the students of applied econometrics who have not had the material in Secs. 1.6 to 1.11, a one-semester course may cover Chaps. 1 to 5, omitting most of the proofs in the starred sections. For the students of theoretical econometrics who have not had the material in Secs. 1.6 to 1.11, a one-semester course may include Chaps. 1 to 5, covering the above sections, and Secs. 2.13, 3.2, and 3.3 *after* the nonstarred sections of Chaps. 1 to 3 have been studied. Less time will be devoted to the empirical studies. Students having had one semester of statistical methods including the materials of Secs. 1.6 to 1.11 may take up additional topics from Chaps. 6, 7, or 9. Chapters 6 to 12 can serve as a text for a second semester in econometrics.

I am indebted to the students at Princeton who helped me learn the material in this text, some using drafts of various chapters and offering comments, including George Mailath, Loretta Mester, and In-Koo Cho, who have also helped prepare the index. Thanks are due to Takeshi Amemiya and John B. Taylor who read Chap. 11, and to David Brownstone who read Chap. 8, and to Adrian

Pagan and Gary Skoog who read many chapters and provided useful comments. It is impossible to express sufficient gratitude to my colleague Richard Quandt, who read the entire manuscript and made many useful suggestions for improvement. Pia Ellen has typed drafts of the manuscript with remarkable efficiency and good spirit. Without her help, the book would not have been completed. I would like to thank all the publishers for granting permission to reprint material written by me and published by them. References are given in the text. Results of research supported by the National Science Foundation through several grants find their way into many sections in this book, especially in Chaps. 7 and 9 to 12.

Gregory C. Chow

CONTENTS

	Preface	xi
Chapter 1	Simple Linear Regression	1
1.1	What Is Econometrics?	1
1.2	Model of Simple Linear Regression	3
1.3	Point Estimation	4
1.4	Testing Hypotheses and Interval Estimation	6
1.5	Use of Matrix Notation	7
*1.6	The Multivariate Normal Distribution and Two Regressions	8
1.7	Errors in Observations	14
*1.8	Convergence in Probability	16
*1.9	Central Limit Theorems	18
*1.10	The Cramer-Rao Inequality	22
*1.11	Asymptotic Distribution of Maximum-Likelihood Estimators	25
1.12	Estimating the Quantity of, and Demand for, Computers	27
Chapter 2	Multiple Linear Regression	38
2.1	The Model of Multiple Linear Regression	38
2.2	Least-Squares Estimation for β and σ^2	39
*2.3	Geometric Interpretation of the Least-Squares Regression	41
2.4	Testing Hypotheses about β and σ^2	43
2.5	Demand for Automobiles: Theory	47
2.6	Demand for Automobiles: Statistical Findings, 1921–1953	53
2.7	Testing the General Linear Hypothesis	58
2.8	Testing Equality between Sets of Regression Coefficients	60
2.9	Forecasting	62
2.10	Testing the Stability of Automobile Demand Functions	64

* Starred sections may be omitted without loss of continuity; see Preface.

2.11	Use of Statistical Demand Functions for Long-Run Forecasting	65
2.12	Partial Correlation Coefficients	69
*2.13	Asymptotic Distribution of the Least-Squares Estimator b	72
Chapter 3	Topics in Regression Analysis	77
3.1	Method of Generalized Least Squares	77
*3.2	Asymptotic Distribution of the GLS Estimator	81
3.3	Analysis of Regression Residuals	84
3.4	Robust Estimators	88
*3.5	Bayesian Estimation	90
3.6	Non-Bayesian Use of Extraneous Information	96
3.7	Multicollinearity without Other Information	98
3.8	Distributed Lags	102
3.9	Errors in Observations	105
3.10	Method of Instrumental Variables	107
Chapter 4	Simultaneous Equations: Model and Identification	111
4.1	Model of Linear Simultaneous Stochastic Equations	111
4.2	Two Problems Associated with the Simultaneous-Equation Model	114
4.3	Conditions for Identifying a Structural Equation	117
*4.4	Identification of a Set of Structural Parameters	122
4.5	Formulation of a Macroeconometric Model	126
4.6	Some Statistical Considerations for the Model	132
4.7	Empirical Results from the Model	137
4.8	Goodness of Fit and Forecasting Value of the Model	141
4.9	Relative Importance of Various Factors in Income Determination	144
4.10	Final Form of a Linear Dynamic Model	145
Chapter 5	Estimation of Linear Simultaneous Equations	153
5.1	Method of Two-Stage Least Squares	153
*5.2	Method of Limited-Information Maximum Likelihood	157
5.3	The k -Class Estimator	164
5.4	Method of Three-Stage Least Squares	167
*5.5	Method of Full-Information Maximum Likelihood	170
5.6	Method of Instrumental Variables	175
*5.7	Treatment of Identities and Linear Restrictions	177
*5.8	FIML with Autoregressive Residuals	178
5.9	Choice of Estimators	180
5.10	Estimation of the Reduced Form	182
Chapter 6	Time-Series Analysis	188
6.1	Time-Series Models	188
6.2	Dynamic Properties of Time Series	190
6.3	Autocovariance Matrix of a Linear Model	192
6.4	Spectral-Density Matrix of a Linear Model	196
*6.5	Decomposition of Time Series into Periodic Components	200

6.6	Note on the Estimation of Spectral Densities	204
*6.7	Estimation of ARMA Models	207
6.8	Box-Jenkins Techniques	211
*6.9	Definition and Tests of Causality	212
Chapter 7	Nonlinear Models	220
7.1	Introduction	220
7.2	Method of GLS or Minimum Distance	222
7.3	Nonlinear Regression	228
7.4	Method of Maximum Likelihood	230
7.5	Numerical Methods of Maximization	232
*7.6	FIML for Nonlinear Simultaneous Equations	235
7.7	Method of Instrumental Variables	240
*7.8	Nonlinear Two- and Three-Stage Least Squares	243
*7.9	Models of Markets in Disequilibrium	244
7.10	Dynamic Properties of Nonlinear Simultaneous Equations	248
Chapter 8	Discrete and Limited Dependent Variables	253
8.1	Introduction	253
8.2	Probit Analysis	254
8.3	Logit Analysis	255
*8.4	Utility Theory for Discrete-Choice Models	257
8.5	Maximum-Likelihood Estimation of Multinomial Logit Models	260
*8.6	Nested Logit Models	263
8.7	Limited Dependent Variables	265
*8.8	The E-M Algorithm	268
8.9	Truncated Sample	271
Chapter 9	Criteria for Model Selection	277
9.1	Introduction	277
9.2	A Method for Selecting Nonnested Regression Models	278
9.3	Some Tests of Nonnested Hypotheses	284
9.4	Lagrangian Multiplier and Related Tests	286
9.5	The C_p Criterion	291
9.6	The Information Criterion	293
9.7	The Posterior-Probability Criterion	300
9.8	Comparison of the Posterior-Probability and Information Criteria	302
*9.9	Estimation of the Information Criterion for Simultaneous-Equation Models	305
9.10	Should a Linear Econometric Model Be Decomposed or Aggregated?	309
9.11	Tests and Analysis of Model Specifications	313
Chapter 10	Models of Time-Varying Coefficients	320
10.1	Introduction	320
10.2	Derivation of β_{it} by Recursive Regression of β_t on y_1, \dots, y_s	321

10.3	Derivations of $\beta_{t s}$ by Regression of y_1, \dots, y_s on x_1, \dots, x_s	326
10.4	Maximum-Likelihood Estimation of σ^2 , V , and M	327
*10.5	System of Linear Regressions with Time-Varying Coefficients	330
*10.6	System of Linear Simultaneous Equations	333
*10.7	System of Nonlinear Simultaneous Equations	337
10.8	Model with Stationary Coefficients	338
*10.9	Identifiability of Parameters	340
10.10	Testing Constancy of Regression Coefficients	342
10.11	The Estimation of Seasonal Components in Economic Time Series	345
Chapter 11	Models under Rational Expectations	351
11.1	The Assumption of Rational Expectations	351
11.2	The Problem of Multiple Solutions	353
11.3	Solution to Linear Expectations Models	356
*11.4	The Solution of Blanchard and Kahn	362
11.5	Estimation of Linear Models without Expectations of Future Variables	364
11.6	Estimation of Linear Models with Future Expectations	366
Chapter 12	Models of Optimizing Agents	373
12.1	Introduction and Preview	373
12.2	Deriving an Optimal Feedback Control Equation	377
*12.3	Method of Maximum Likelihood	380
12.4	Method of Two-Stage Least Squares	385
12.5	Two Economic Examples	386
12.6	Alternative Derivation of the Optimal Rule	390
12.7	Explicit Solution for the Optimal Rule	392
12.8	The Assumptions of Optimization Models	394
12.9	Model of a Dynamic Game	395
*12.10	Policy Evaluation and Optimization under Rational Expectations	396
*12.11	Estimation of a Dynamic Game Model with a Dominant Player	400
*12.12	Estimation of a Dynamic Game Model under Nash Equilibrium	401
	Tables	405
	Index	419

SIMPLE LINEAR REGRESSION

1.1 WHAT IS ECONOMETRICS?

Econometrics is the art and science of using statistical methods for the measurement of economic relations. In the practice of econometrics, economic theory, institutional information, and other assumptions are relied upon to formulate a statistical model, or a set of statistical hypotheses, to explain the phenomena in question. Econometric methods are used to estimate the parameters of the model, to test hypotheses concerning them, and to generate forecasts from the model. The formulation of an econometric model is an art, just as using knowledge of architecture to design a building is an art. A good econometrician uses sound judgment to bring the relevant knowledge in economics to bear in formulating a useful model. The most important variables are selected while the nonessential ones are discarded. The crucial relationships are formulated and incorporated in the model. Care is taken to ensure that the statistical data used actually correspond to the variables to be measured according to theoretical considerations.

Given a set of requirements in the construction of a building, two good architects will come up with two different designs. Both may serve the purposes well. Similarly, given the same objectives, two econometricians are not likely to come up with two identical models, but both may capture the essential elements of the problem sufficiently to be useful. Good econometric models, like good architectural designs, can serve as prototypes to be followed in future investigations. The art of formulating a good econometric model is difficult to learn. One needs a solid command of the tools of economic analysis and sound judgment to select the essential variables of the problem. In a Walrasian system of general equilibrium all economic variables are related, but only a subset of variables will be selected in a particular investigation. One can read the best econo-

metric studies to see how they were done, in the same way that an architect reads the best designs or an artist studies masterpieces of art. Finally, one can practice building one's own models and thus learn by doing.

This book is devoted mainly to an easier aspect of the practice of econometrics, namely, the study of some of the statistical methods most useful for drawing inferences from an econometric model. Knowledge of econometric methods is essential, though by no means sufficient, for the construction of good econometric models, as we have just pointed out. Our presentation presumes that the reader is familiar with the basic ideas of statistical inference. To review some of the basic ideas and to indicate the level of preparation suitable for reading this book, we shall introduce the required tools of statistics in this chapter and apply them to the model of simple linear regression as a warmup exercise. Most of the techniques presented later in this book can be viewed as generalizations, extensions, or modifications of this simple model.

At the beginning of an econometric investigation, an econometrician needs to know clearly what economic phenomena are to be explained, what important factors will contribute to the explanation, how these factors should be measured, what quantitative relationships exist, how such relations can be estimated or tested, and what conclusions can be drawn from the investigation. Consider the example of studying the demand for apples in the United States. The annual consumption of apples through time may be selected as the phenomenon to be explained. The important factors explaining the demand for apples may be the price of apples and income of the consumer. One may choose the per capita annual consumption of apples (in pounds) as the dependent variable. The price of apples per pound, deflated by a consumer price index, and per capita disposable income, also deflated by a consumer price index, may serve as the explanatory variables. *The logarithm of per capita apple consumption may be assumed to be a linear function of the log of relative price, the log of per capita real disposable income, and a normally distributed random disturbance summarizing the combined effects of omitted factors.* Time-series data may be used to estimate this linear relation. The coefficients of log price and log income will be interpreted respectively as the price and income elasticities of demand for apples. One objective of the investigation may be to estimate these demand elasticities. Another may be to test the hypotheses that both elasticities are less than 1 in absolute value. A third may be to use the estimated relation to forecast the demand for apples 5 years from now.

In the above illustration, many issues will have to be resolved by the econometrician. For example, is deflation of total apple consumption by total population in the United States sufficient to account for the effects of demographic factors on apple consumption? Is the age distribution relevant? Is mean income sufficient to explain apple consumption without taking the distribution of income into account? Is the relationship approximately linear in the logarithms of the variables? Does one have to account for the effects of lagged incomes on consumption? Is it useful to combine cross-section data with time-series data to estimate the price and income elasticities? Does the coefficient of log price in the postulated linear relation measure solely the price elasticity of demand, rather

than the elasticity of supply? Is it necessary to formulate a supply equation for apples and estimate a system of two equations simultaneously?

The study of demand for apples is an easy econometric problem, and yet all these issues and others must be resolved. Consider the more difficult problem of studying the demand for computers in the United States. Just the measurement of the dependent variable requires some serious thought. Since there are large and small computers, or computers with different computing powers and capabilities, how is the total quantity of computers to be measured? Since the quality of the computers, however measured, has improved so rapidly through time, how should the effect of this rapid technological change be incorporated in a study of the demand for computers? Related to the measurement of the quantity of computers is the problem of measuring the price per unit of computers. If the quantity of computers can be measured appropriately, one will know what 1 unit of computers is and the price per unit can be determined accordingly. These problems are only illustrative of the issues an econometrician faces in the formulation of an econometric model. Other issues have been stated at the beginning of the paragraph before last. The problem of the demand for computers is considered in Sec. 1.12. The issues in formulating an econometric model cannot be systematically discussed in this book but will be illustrated by case studies cited in Chaps. 1, 2, and 4. Before embarking on a statistical analysis, the econometrician should have resolved all these issues and planned ahead, deciding how conclusions should be drawn before performing the statistical computations. Having made these introductory remarks, we begin by describing the model of simple linear regression.

1.2 MODEL OF SIMPLE LINEAR REGRESSION

Let y_i denote the i th observation of the dependent variable and x_i denote the associated explanatory variable. A simple linear regression model can be written as

$$y_i = \alpha + \beta x_i + \epsilon_i \quad i = 1, \dots, n \quad (1)$$

where the residual or disturbance terms ϵ_i are assumed to be normal and independent, each having mean zero and variance σ^2 . The observations on x_i are treated as fixed numbers, except when specified otherwise, as in Secs. 1.6, 1.7, and 3.8. That is, a probability distribution for x_i is not postulated. The parameters of this model are α , β , and σ^2 . The mean of the random variable y_i is $\alpha + \beta x_i$, as we observe by taking the mathematical expectations (or means) of both sides of (1). The mean is thus a linear function of x ; the mean is called a *regression function*. As an example, consider n families selected from Princeton, New Jersey. Let y_i denote the logarithm of the quantity of apples consumed by the i th family during a given period and x_i denote the logarithm of its income during the same period. The mean of y_i can be assumed to be a linear function of x_i . Given x_i , that is, for families having log income equal to x_i , the distribution of log apple consumption is assumed to be normal, having mean $\alpha + \beta x_i$ and standard deviation σ .

In classical statistics the two main problems are to estimate the unknown

parameters, α , β , and σ^2 in this case, and to test hypotheses concerning them. These topics will be discussed in turn.

1.3 POINT ESTIMATION

Let point estimates of the parameters α and β be denoted respectively by a and b . Other methods will be presented in due course, but first we discuss a popular method for finding these estimates, namely, the method of least squares. In this method one chooses the values of a and b which minimize the sum of squares of the deviations of y_i from the estimated regression line $a + bx_i$, that is,

$$\sum_{i=1}^n [y_i - (a + bx_i)]^2$$

Setting to zero the derivative of this sum with respect to a , we obtain

$$a = \frac{\sum_{i=1}^n y_i - b \sum_{i=1}^n x_i}{n} \equiv \bar{y} - b\bar{x} \quad (2)$$

Substituting $\bar{y} - b\bar{x}$ for a in this sum and setting to zero its derivative with respect to b , we obtain

$$\begin{aligned} b &= \left[\sum_i x_i(x_i - \bar{x}) \right]^{-1} \sum_i x_i(y_i - \bar{y}) \\ &= \left[\sum_i (x_i - \bar{x})^2 \right]^{-1} \sum_i (x_i - \bar{x})(y_i - \bar{y}) \end{aligned} \quad (3)$$

where we have used $\sum_i \bar{x}(x_i - \bar{x}) = 0$ and $\sum_i \bar{x}(y_i - \bar{y}) = 0$. Equations (2) and (3) are the *normal equations* for computing the estimates a and b .

To decide whether these formulas or *estimators* are good we consider their sampling distributions. As seen from (2) and (3), the estimators a and b are functions of the random variables y_1, \dots, y_n and therefore random variables themselves. The term *estimates* refers to the particular numbers representing a and b which we obtain by using a particular sample (y_1, \dots, y_n) . The term *estimators* refers to possible numbers a and b which we could obtain by (hypothetically) drawing repeated samples of (y_1, \dots, y_n) given the values of x_1, \dots, x_n . An estimator is a random variable; an estimate is a number which is computed by using the data for one particular sample. To evaluate the mean and variance of the distribution of b or of the random variable b we substitute $\alpha + \beta x_i + \epsilon_i$ for y_i in (3) and write b as a function of ϵ_i

$$\begin{aligned} b &= \left[\sum_i (x_i - \bar{x})^2 \right]^{-1} \left[\sum_i (x_i - \bar{x}) y_i \right] \\ &= \left[\sum_i (x_i - \bar{x})^2 \right]^{-1} \left[\sum_i (x_i - \bar{x})(\alpha + \beta x_i + \epsilon_i) \right] \end{aligned}$$

$$\begin{aligned}
&= \left[\sum_i (x_i - \bar{x})x_i \right]^{-1} \left[\sum_i (x_i - \bar{x})\beta x_i + \sum_i (x_i - \bar{x})\epsilon_i \right] \\
&= \beta + \left[\sum_i (x_i - \bar{x})^2 \right]^{-1} \left[\sum_i (x_i - \bar{x})\epsilon_i \right] \quad (4)
\end{aligned}$$

The mean of the random variable b is obtained by taking the mathematical expectation of (4) yielding

$$E[b] = E[\beta] + \left[\sum_i (x_i - \bar{x})^2 \right]^{-1} \left[\sum_i (x_i - \bar{x})E\epsilon_i \right] = \beta \quad (5)$$

where we have treated $\sum_i (x_i - \bar{x})^2$ as a constant and taken the expectation of the sum $\sum_i (x_i - \bar{x})\epsilon_i$ as the sum $\sum_i E[(x_i - \bar{x})\epsilon_i]$ of the expectations, each $E\epsilon_i$ being zero by assumption. Since $E[b] = \beta$, b is an *unbiased estimator* of β .

To find the variance of b we use (4) and (5) to evaluate

$$\begin{aligned}
\text{Var } b &\equiv E[(b - Eb)^2] = E \left\{ \left[\sum_i (x_i - \bar{x})^2 \right]^{-1} \left[\sum_i (x_i - \bar{x})\epsilon_i \right] \right\}^2 \\
&= \left[\sum_i (x_i - \bar{x})^2 \right]^{-2} E \left[\sum_i (x_i - \bar{x})\epsilon_i \right]^2 \\
&= \left[\sum_i (x_i - \bar{x})^2 \right]^{-2} \sum_i (x_i - \bar{x})^2 E\epsilon_i^2 = \left[\sum_i (x_i - \bar{x})^2 \right]^{-1} \sigma^2 \quad (6)
\end{aligned}$$

where the third line has used the assumption that $E\epsilon_i\epsilon_j = 0$ for $i \neq j$. Thus the variance of the sampling distribution of b is directly proportional to the variance σ^2 of the regression residuals ϵ_i and inversely proportional to the variance of the explanatory variable. Having more variations in x permits a tighter estimate of b .

Note from the first line of (4) that b is a linear function of y_1, \dots, y_n and is called a *linear estimator*. It is an *unbiased estimator* of β , as shown by (5). In fact, among all linear, unbiased estimators, b has the *smallest variance*, as will be proved in Chap. 2. Thus, the least-squares estimator b is said to be *best linear unbiased*. We leave the evaluation of the mean and variance of the distribution of the estimator a as exercises.

A frequently used point estimate of σ^2 is

$$s^2 = (n - 2)^{-1} \sum_{i=1}^n (y_i - a - bx_i)^2 \quad (7)$$

It will be shown in Chap. 2 that $(n - 2)s^2/\sigma^2$ or $\sum_{i=1}^n (y_i - a - bx_i)^2/\sigma^2$ is distributed as χ^2 with $n - 2$ degrees of freedom. Since a $\chi^2(n - 2)$ distribution has a mean equal to $n - 2$, the mean of s^2/σ^2 is 1, or the mean of s^2 is σ^2 . In other words, s^2 is an unbiased estimator for σ^2 . However, the mean-squared error $E(s^2 - \sigma^2)^2$ can be improved by replacing the denominator $n - 2$ by $n - 1$. This would produce a biased estimator, whose expectation is smaller than σ^2 , but the variance of the estimator would also be reduced, leading to a smaller mean-squared error. The estimator s^2 is often used because it is convenient to use standard tables for the χ^2 distribution for its sampling distribution.

1.4 TESTING HYPOTHESES AND INTERVAL ESTIMATION

From (4) we observe that $b - \beta$ is a linear combination of $\epsilon_1, \dots, \epsilon_n$. Therefore, if the ϵ_i are normal, the linear combination $b - \beta$ will also be normal. Even if ϵ_i are not normal but have a finite variance, the linear combination $b - \beta$ will have a distribution which approaches the normal distribution as the sample size n increases, according to the central limit theorem, discussed in Sec. 1.9. From (6) we obtain the variances of b . Hence

$$\frac{b - \beta}{\left[\sum_i (x_i - \bar{x})^2\right]^{-1/2}\sigma} \quad (8)$$

has a standard normal distribution. If σ^2 is known, the statistic (8) can be used to test hypotheses concerning β . Let the null hypothesis be $\beta = \beta_0$ and the alternative hypothesis be $\beta \neq \beta_0$. Let the level of significance be 5 percent. Under the null hypothesis, the probability that the statistic (8) with $\beta = \beta_0$ will be larger than 1.96 or smaller than -1.96 is .05. If this statistic as computed from our sample turns out to exceed 1.96 in absolute value, the null hypothesis $\beta = \beta_0$ will be rejected. If the alternative hypothesis is one-sided, say $\beta < \beta_0$, the null hypothesis will be rejected when the statistic (8) is smaller than -1.64 since the probability for a standard normal random variable to be smaller than this value is .05.

To construct a symmetric interval estimate or confidence interval for β with a confidence coefficient of .95 we use inequalities

$$-1.96 < \frac{b - \beta}{\left[\sum_i (x_i - \bar{x})^2\right]^{-1/2}\sigma} < 1.96$$

implying

$$b + 1.96 \left[\sum_i (x_i - \bar{x})^2\right]^{-1/2} \sigma > \beta > b - 1.96 \left[\sum_i (x_i - \bar{x})^2\right]^{-1/2} \sigma \quad (9)$$

To test the null hypothesis that $\sigma^2 = \sigma_0^2$ we use the fact that under the null hypothesis $(n - 2)s^2/\sigma_0^2$ has a $\chi^2(n - 2)$ distribution. If this statistic exceeds an upper critical value C_U or falls below a lower critical value C_L according to the $\chi^2(n - 2)$ distribution, the null hypothesis will be rejected. Similarly a confidence interval for σ^2 can be constructed by using

$$C_L < \frac{(n - 2)s^2}{\sigma^2} < C_U$$

implying

$$\frac{(n - 2)s^2}{C_L} > \sigma^2 > \frac{(n - 2)s^2}{C_U}$$

To test the null hypothesis $\beta = \beta_0$ when σ^2 is unknown we use the statistic

$$\frac{b - \beta_0}{\left[\sum_i (x_i - \bar{x})^2\right]^{-1/2}s} \quad (10)$$

which is the ratio of (8) to s/σ . The numerator (8) is standard normal. The denominator s/σ is the square root of a $\chi^2(n-2)$ variable divided by the degrees of freedom $(n-2)$. As will be shown in Chap. 2, the numerator and the denominator are statistically independent. Therefore, the ratio (10) will have Student's t distribution with $n-2$ degrees of freedom. Hypothesis testing and interval estimation for β can be performed using the $t(n-2)$ distribution for the statistic (10). For example, the .025 upper-tail critical values for $t(10)$ and $t(20)$ are respectively 2.228 and 2.086, compared with 1.960 for the standard normal statistic. A $\chi^2(m)$ variable is the sum of squares of m independent standard normal variables. If $\chi^2(m)$ and $\chi^2(n)$ are independent, the ratio of $\chi^2(m)/m$ to $\chi^2(n)/n$ has an $F(m, n)$ distribution. An $F(1, n)$ variable is the square of a $t(n)$ variable. Tables for the normal, χ^2 , t , and F distributions are provided at the end of this book.

1.5 USE OF MATRIX NOTATION

In order to study the model of multiple linear regression, where the mean of y is a linear function of several explanatory variables, it is convenient to use matrix notation. Let the model (1) be written as

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i \quad i = 1, \dots, n \quad (11)$$

where we write the original α as β_1 and the original β as β_2 , and let the i th observation of the first independent variable x_{i1} be equal to 1 identically so that its coefficient β_1 is the intercept. Denote by y , x_1 , x_2 , and ϵ , respectively, the column vectors of n observations on the dependent variable, the first and the second explanatory variables, and the random disturbance, namely,

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad x_1 = \begin{bmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{bmatrix} \quad x_2 = \begin{bmatrix} x_{12} \\ x_{22} \\ \vdots \\ x_{n2} \end{bmatrix} \quad \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

We can write the n observations of the model (11) as

$$y = \beta_1 x_1 + \beta_2 x_2 + \epsilon \quad (12)$$

The vector y equals the sum of three vectors, $\beta_1 x_1$, $\beta_2 x_2$, and ϵ . When a scalar β_1 multiplies a vector x_1 , each element of x_1 is multiplied by the scalar to form a vector $\beta_1 x_1$ as the product. When the three vectors $\beta_1 x_1$, $\beta_2 x_2$, and ϵ are added, their corresponding elements are added to form the elements of the sum, which in our case equals the vector y . The equality sign in (12) signifies that each element of y equals the corresponding element of the vector on the right, obtained as the sum of the three vectors.

The inner product of a row vector and a column vector is the sum of the products of their corresponding elements. For example, let x'_1 be the row vector $[x_{11} \ x_{21} \ \cdots \ x_{n1}]$ obtained by taking the transpose (indicated by a prime) of the column vector x_1 , and let x_2 be the column vector as defined above. The