## Diego Escobari

The University of Texas - Pan American

# Introduction to Econometrics

(preliminary class notes)

## ECON 3341

February 16, 2012

## Contents

1	Ran	dom Vai	riables, Sampling and Estimation	1
	1.1	Introdu	action	1
	1.2	Probab	vilities	1
		1.2.1	Events	1
		1.2.2	Probability postulates	2
	1.3	Discret	te random variables and expectations	3
		1.3.1	Discrete random variables	3
		1.3.2	Expected value of random variables	4
		1.3.3	Expected value rules	5
		1.3.4	Variance of a discrete random variable	5
		1.3.5	Probability density	6
	1.4	Contin	uous random variables	7
		1.4.1	Probability density	7
		1.4.2	Normal distribution	8
		1.4.3	Expected value and variance of a continuous random variable	9
	1.5	Covari	ance and correlation	9
		1.5.1	Covariance	9
		1.5.2	Correlation	9
	1.6	Sampli	ing and estimators	10
		1.6.1	Sampling	10
		1.6.2	Estimators	11
	1.7	Unbias	sedness and efficiency	11
		1.7.1	Unbiasedness	11
		1.7.2	Efficiency	12
		1.7.3	Unbiasedness versus efficiency	13
	1.8	Estima	tors for the variance, covariance, and correlation	13
	1.9	Asymp	ototic properties of estimators	14
		1.9.1	Consistency	15
		1.9.2	Central limit theorem	15

## Contents

2	Sim	ple Linear Regression	17
	2.1	Simple linear model	17
	2.2	Least squares regression	17
	2.3	Interpretation of the regression coefficients	19
	2.4	Goodness of fit	20
3	Prop	perties and Hypothesis Testing	23
	3.1	Types of data	23
	3.2	Assumptions of the model	23
	3.3	Unbiasedness of the coefficients	25
	3.4	Precision of the coefficients	25
	3.5	The Gauss-Markov theorem	26
	3.6	Hypotheses testing	26
		3.6.1 Formulation of the null hypothesis	26
		3.6.2 <i>t</i> -tests	27
		3.6.3 Confidence intervals	28
		3.6.4 <i>F</i> test	29
	3.7	Computer output	31
			-
4	Mul	tiple Regression Analysis	33
	4.1	Interpretation of the coefficients	33
	4.2	Ordinary Least Squares	34
	4.3	Assumptions	34
	4.4	Properties of the coefficients	35
		4.4.1 Unbiasedness	35
		4.4.2 Efficiency	35
		4.4.3 Precision of the coefficient, $t$ tests, and confidence intervals.	35
	4.5	Regression output in Gretl	36
	4.6	Multicollinearity	37
	4.7	Goodness of fit: $R^2$ and $R^2$	38
	4.8	F tests	38
	4.9	Adjusted $R^2$ , $R^2$	39
5	Trar	sformations of Variables and Interactions	41
	5.1	Basic idea	41
	5.2	Logarithmic transformations	42
	5.3	Quadratic terms	43
	5.4	Interaction terms	45
6	Ana	lysis with Qualitative Information: Dummy Variables	47
v	6.1	Describing qualitative information	47
	6.2	A single dummy independent variable	48
	63	Dummy variables for multiple categories	50
	64	Incorporating ordinal information	51
	6.5	Interactions involving dummy variables	52
	0.5	6 5 1 Allowing for different slopes	52
			55

vi

		6.5.2	Testing for differences in regression functions across groups	54
	6.6	The du	mmy variable trap	57
-	<b>C</b>	:C	A Deservation Variables	50
/	Spec	<i>yıcanor</i> Mədəl	and of Kegression variables	59 50
	7.1	Omitti		39 50
	1.2		The bigs problem	50
		7.2.1	Ine blas problem	39 60
		7.2.2		60
	7 2	1.2.3		60
	1.3	Includi	ng a variable that should not be included	63
	7.4	7.3.1	Example	64
	7.4	Testing	g a linear restriction	65
8	Hete	rosceda	sticity	69
Ū	8 1	Hetero	scedasticity and its implications	69
	8.2	Testing	for heteroscedasticity	69
	0.2	8 2 1	Breusch-Pagan test	69
		822	Breusch Pagan test in Gret	71
		822	White test	71
		824	White test in Gret	72
	83	What t	o do with beteroscedosticity?	72
	0.5	8 2 1	Simple transformation of the variables	72
		0.3.1	Weighted Least Squares	73 72
		0.3.2 0.2.2	Weighted Least Squares in Croth	15
		0.3.3	White's heterogoodosticity consistent standard arrays	74 76
		8.3.4	white's neteroscedasticity-consistent standard errors	70
9	Bina	rv Choi	ce Models	77
	9.1	The lin	ear probability model	77
		9.1.1	The model	77
		9.1.2	The linear probability model in Gret	78
	9.2	Logit a	nalvsis	79
		9.2.1	The logit transformation	79
		922	Logit regression in Gretl	80
	9.3	Probit	analysis	81
	1.0	9.3.1	The probit transformation	81
		9.3.2	Probit regression in Gret	81
		,		
10	Time	Series		83
	10.1	Time S	eries Data	83
	10.2	Time S	eries Regression Models	84
		10.2.1	Static Models	84
		10.2.2	Finite Distributed Lag Models	85
		10.2.3	Autoregressive Model	87
		10.2.4	Moving-Average Models	88
		10.2.5	Autoregressive Moving Average Models	91
		-		

vii

## Chapter 1 Random Variables, Sampling and Estimation

## **1.1 Introduction**

This chapter will cover the most important basic statistical theory you need in order to understand the econometric material that will be coming in the next chapters. The key topics that we will review are the following:

- Descriptive statistics. e.g. mean and variance.
- *Probability*. e.g. events, relative frequency, marginal and conditional probability distributions.
- Random variables, probability distributions, and expectations.
- Sampling. e.g. simple random sampling.
- Estimation. e.g. the distinction between and estimator and an estimate.
- Statistical inference. t and F tests.

## **1.2 Probabilities**

## 1.2.1 Events

**Random experiment**. Process leading to two or more possible outcomes, with uncertainty as to which outcome will occur.

Flip of a coin, toss of a die, a students takes a class and either obtains an A or not.

Sample space. Set of all basic outcomes of a random experiment.

When flipping a coin, S = [head, tail]. When taking a class, S = [A, B, C, D, F, drop]. When tossing a die, S = [1, 2, 3, 4, 5, 6]. No two outcomes can occur simultaneously.

Event. Subset of basic outcomes in the sample space.

Event E<sub>1</sub>: "Pass the class" then the subset of basic outcomes is A, B, C.

**Intersection of event**. When two events  $E_1$  and  $E_2$  have some basic outcomes in common. It is denoted by  $E_1 \cap E_2$ .

Event  $E_1$ : Individuals with college degree. Event  $E_2$ : Individuals who are married.  $E_1 \cap E_2$ : Individuals who have college degree and are married.

Joint probability. Probability that the intersection occurs.

**Mutually exclusive events**.  $E_1$  and  $E_2$  are mutually exclusive if  $E_1 \cap E_2$  is empty.

**Union of events**. Denoted by  $E_1 \cup E_2$ . At least one of these events occurs. Either  $E_1$ ,  $E_2$ , or both.

**Complement**. The complement of E is denoted by  $\overline{E}$  and it is the set of basic outcomes of a random experiment that belongs to S, but not to  $E_1$ .

 $E_1$  is the complement of  $\overline{E}_1$ Event  $E_2$ : Individuals who are married.  $E_1$  and  $\overline{E}$  are mutually exclusive events.

## **1.2.2** Probability postulates

Given a random experiment, we want to determine the probability that a particular event will occur. A probability is a measure from 0 to 1.

2

1.3 Discrete random variables and expectations

 $0 \rightarrow$  the event will not occur.  $1 \rightarrow$  the event is certain.

When the outcomes are equally likely to occur, the probability of an event E is:

$$\begin{split} P(E) &= N_E / N \\ N_E: \text{Number of outcomes in event E.} \\ N: \text{ Total number of outcomes in the sample space S.} \end{split}$$

Example 1: Flip of a coin, Event E is "head" then P(E) = 1/2.  $N_E = 1$  and N = 2.

Example 2: Event E is "winning the lottery" then if there are 1000 lottery tickets and you bought, 2 P(E) = 2/1000 = 0.002.

#### Some probability rules

$$\begin{split} P(E\cup\bar{E}) &= P(E) + P(\bar{E}) = 1.\\ P(\bar{E}) &= 1 - P(E). \end{split}$$

Conditional probability

 $P(E_1 | E_2)$ : Probability that  $E_1$  occurs, given that  $E_2$  has already occurred.  $P(E_1 | E_2) = P(E_1 \cap E_2) / P(E_2)$  given that  $P(E_2) > 0$ .

Addition rule

 $P(E_1 \cup E_2) = P(E_1) + P(E_2) - P(E_1 \cap E_2).$ 

Statistically independent events

 $\begin{aligned} P(E_1 \cap E_2) &= P(E_1)P(E_2). \\ P(E_1 \mid E_2) &= P(E_1)P(E_2) / P(E_2) = P(E_1). \end{aligned}$ 

## 1.3 Discrete random variables and expectations

## 1.3.1 Discrete random variables

**Random variable**. Variable that takes numerical values determined by the outcome of a random experiment.

Examples: Hourly wage, GDP, inflation, the number when tossing a die. Notation: Random variable *X* can take *n* possible values  $x_1, x_2, \dots x_n$ .

**Discrete random variable**. A random variable that takes a countable number of values.

Examples: Number of years of education.

**Continuous random variable**. A random variable that can take any value on an internal.

Examples: Wage, GDP, exact weight.

Consider tossing two dies (green and red). This will yield 36 possible outcomes because the green can take 6 possible values and the red can take also 6 values,  $6 \times 6 = 36$ . Let's define the random variable *X* to be the sum of two dice. Therefore *X* can take 11 possible values, from 2 to 12. This information is summarized in the following tables.

Table 1.1 Outcomes with two dies

red / green	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8
3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12

Table 1.2 Frequencies and probability distributions

Value of X	2	3	4	5	6	7	8	9	10	11	12
Frequency	1	2	3	4	5	6	5	4	3	2	1
Probability (p)	1/36	2/36	3/36	4/36	5/36	6/36	5/36	4/36	3/36	2/36	1/36

## 1.3.2 Expected value of random variables

Let E(X) be the expected value of the random variable X. The expected value of a discrete random variable is the weighted average of all its possible values, taking the probability of each outcome as its weight. Random variable X can take n particular values  $x_1, x_2, \ldots, x_n$  and the probability of  $x_i$  is given by  $p_i$ . Then we have that the expected value is given by:

$$E(X) = x_1 p_1 + x_2 p_2 + \dots + x_n p_n = \sum_{i=1}^n x_i p_i.$$
 (1.1)

We can also write the expected value as:  $E(X) = \mu_X$ . For the previous example we can calculate that the expected value as:

4

1.3 Discrete random variables and expectations

$$E(X) = 2 \cdot \frac{1}{36} + 3 \cdot \frac{2}{36} + \dots + \frac{12 \cdot \frac{1}{36}}{6} = \frac{252}{36} = 7$$
(1.2)

X	р	Xp
2	1/36	2/36
3	2/36	6/36
4	3/36	12/36
5	4/36	20/36
6	5/36	30/36
7	6/36	42/36
8	5/36	40/36
9	4/36	36/36
10	3/36	30/36
11	2/36	22/36
12	1/36	12/36
Total	$E(X) = \sum_{i=1}^{n} x_i p_x$	252/36 = 7

 Table 1.3 Expected value of X, two dice example

## 1.3.3 Expected value rules

Let X, Y, and Z denote three random variables, and let b,  $b_1$ , and  $b_2$  be arbitrary constants. Then,

$$E(X + Y + Z) = E(X) + E(Y) + E(Z)$$
(1.3)

$$E(bX) = bE(X)$$
 for a constant  $b$  (1.4)

$$E(b) = b \tag{1.5}$$

For the example where  $Y = b_1 + b_2 X$ ,  $b_1$  and  $b_2$  are constants we want to calculate E(X).

$$E(Y) = E(b_1 + b_2 X)$$
(1.6)  
=  $E(b_1) + E(b_2 X)$   
=  $b_1 + b_2 E(X)$ 

## 1.3.4 Variance of a discrete random variable

Let var(X) be the variance of the random variable *X*. var(X) is a useful measure of the dispersion of its probability distribution. It is defined as the expected value of the square of the difference between *X* and its mean. That is,  $E[(X - \mu_X)^2]$ , where

#### 1 Random Variables, Sampling and Estimation

 $\mu_X$  is the population mean of X.

$$\operatorname{var}(X) = \sigma_X^2 = E[(X - \mu_X)^2] \tag{1.7}$$

$$= (x_1 - \mu_X)^2 p_1 + (x_2 - \mu_X)^2 p_2 + \dots + (x_n - \mu_X)^2 p_n \qquad (1.8)$$
$$= \sum_{i=1}^n (x_i - \mu_X)^2 p_i$$

Taking the square root of the variance  $(\sigma_X^2)$  one can obtain the standard deviation,  $\sigma_X$ . The standard deviation also serves as a measure of dispersion of the probability distribution. A useful way to write the variance is:

$$\sigma_X^2 = E(X^2) - \mu_X^2. \tag{1.9}$$

From the previous example of tossing two dies, we have that the population variance can be calculated as follows:

 Table 1.4 Population variance, X from the two dice example

X	р	$X - \mu_X$	$(X-\mu_X)^2$	$(X-\mu_X)^2 p$
2	1/36	-5	25	0.69
3	2/36	-4	16	0.89
4	3/36	-3	9	0.75
5	4/36	-2	4	0.44
6	5/36	-1	1	0.14
7	6/36	0	0	0.00
8	5/36	1	1	0.14
9	4/36	2	4	0.44
10	3/36	3	9	0.75
11	2/36	4	16	0.89
12	1/36	5	25	0.69
Total				5.83

## 1.3.5 Probability density

Because discrete random variables, by definition, can only take a finite number of values, they are easy to summarize graphically. The probability distribution is the graph that links all the values that a random variable can take with its corresponding probabilities. For the two dice example above, see Figure 1.1.



Fig. 1.1 Discrete probabilities, X from the two dice example

## 1.4 Continuous random variables

### 1.4.1 Probability density

Continuous random variables can take any value on an interval. This means that it can take an infinite number of different values, hence it is not possible to obtain a graph like the one presented in Figure 1.1 for a continuous random variable. Instead, we will define the probability of a random variable lying within a given interval. For example, the probability that the height of an individual is between 5.5 and 6 feet. This is depicted in Figure 1.2 as the shaded area below the probability density curve for the values of X between 5.5 and 6. The probability of the random variable X written as a function of the random variable is known as the probability density function. We can write this ones as f(X). Then, if we use a little math we can easily find the area under the curve. Recall that the are under a curve can be obtained by taking the integral.

**Probability density function**. Is a function that describes the relative likelihood for a random variable to occur at a given point.

$$\int_{5.5}^{6} f(X) = 0.18$$
(1.10)  
$$\int_{0}^{\infty} f(X) = 1$$

#### 1 Random Variables, Sampling and Estimation



Fig. 1.2 Continuous probabilities, X from the height example

The first line in the equation above just calculates the integral under the curve f(X) between the points 5.5 and 6. The second line shows that the whole area under the curve presented in Figure 1.2 is equal to one. This is for the same reason why the summation of all the bars in Figure 1.1 are also equal to one; the total probability is always equal to one.

## 1.4.2 Normal distribution

The normal distribution is the most widely known continuous probability distribution. The graph associated with its probability density function has a bell-shape and its is known as the Gaussian function or bell curve. Its probability density function is given by:

$$f(X) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu^2)}{2\sigma^2}}$$
(1.11)

where  $\mu$  is the mean and  $\sigma^2$  is the variance. Figure 1.2 is an example of this distribution.

1.5 Covariance and correlation

# 1.4.3 Expected value and variance of a continuous random variable

The basic difference between a discrete and a continuous random variable is that the second can take on infinite possible values, hence the summations signs that are used to calculate the expected value and the variance of a discrete random variable cannot be used for a continuous random variable. Instead, we use integral signs. For the expected value we have:

$$E(X) = \int X f(X) dX \tag{1.12}$$

where the integration is performed over the interval for which f(X) is defined. For the variance we have:

$$\sigma_X^2 = E[(X - \mu_X)^2] = \int (X - \mu_X)^2 f(X) dX$$
(1.13)

## 1.5 Covariance and correlation

## 1.5.1 Covariance

When dealing with two variables, the first question you want to answer is whether these variables move together or whether they move in opposite directions. The covariance will help us answer that question. For two random variables X and Y, the covariance is defined as:

$$\operatorname{cov}(X,Y) = \sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)]$$
(1.14)

where  $\mu_X$  and  $\mu_Y$  are the population means of *X* and *Y*, respectively. When to random variables are independent, their covariance is equal to zero. When  $\sigma_{XY} > 0$  we say that the variables move together. When  $\sigma_{XY} < 0$  they move in opposite directions.

## 1.5.2 Correlation

One concern when using the cov(X, Y) as a measure of association is that the result is measured in the units of X times the units of Y. The correlation coefficient, that is dimensionless, overcomes this difficulty. For variables X and Y the correlation coefficient is defined as: 1 Random Variables, Sampling and Estimation

$$\operatorname{corr}(X,Y) = \rho_{YX} = \frac{\sigma_{YX}}{\sqrt{\sigma_X^2 \sigma_Y^2}}$$
(1.15)

The correlation coefficient is a number between -1 and 1. When it is positive, we say that there is a positive correlation between *X* and *Y* and that these two variables move in the same direction. When it is negative, we say that they move in opposite directions.

## 1.6 Sampling and estimators

Notice that in the two dice example we know the population characteristics, that is, the probability distribution. From this probability distribution it is easy to obtain the population mean an variance. However, what happens most of the time is that we need to rely on a data set to get estimates of the population parameters (e.g the mean and the variance). In that case the estimates of the population parameters are obtained using estimators, and the sample needs to have certain characteristics. The estimators and the sampling are the subject of this section.

## 1.6.1 Sampling

The most common way to obtain a sample from the population is through simple random sampling.

**Simple random sampling**. It is a procedure to obtain a sample from the population, where each of the observations is chosen randomly and entirely by chance. This means that each observation in the population has the same probability of being chosen.

Once the sample of the random variable *X* has be generated, each of the *n* observations can be denoted by  $\{x_1, x_2, \dots, x_n\}$ .<sup>1</sup>

<sup>&</sup>lt;sup>1</sup> The textbook Dougherty (2007) makes the distinction between the specific values of the random variable X before and after they are known, and emphasizes this distinction by using uppercase and lowercase letter. This distinction is useful only in some cases and that is why most textbooks do not make this distinction. We will follow emphasize the distinction and we will use only lowercase letters.

1.7 Unbiasedness and efficiency

## 1.6.2 Estimators

**Estimator**. It is a general rule (mathematical formula) for estimating an unknown population parameter given a sample of data.

For example, an estimator for the population mean is the sample mean:

$$\bar{X} = \frac{1}{n}(x_1 + x_2 + \dots + x_n) = \frac{1}{n}\sum_{i=1}^n x_i.$$
(1.16)

An interesting feature of this estimator is that the variance of  $\bar{X}$  is 1/n times the variance of X. The derivation is the following:

$$\sigma_{\bar{X}}^2 = \operatorname{var}(\bar{X}) \tag{1.17}$$

$$\sigma_{\bar{X}}^2 = \operatorname{var}\{\frac{1}{n}(x_1 + x_2 + \dots + x_n)\}$$
(1.18)

$$\sigma_{\bar{X}}^2 = \frac{1}{n^2} \operatorname{var}\{\frac{1}{n}(x_1 + x_2 + \dots + x_n)\}$$
(1.19)

$$\sigma_{\bar{X}}^2 = \frac{1}{n^2} \{ \operatorname{var}(x_1) + \operatorname{var}(x_2) + \dots + \operatorname{var}(x_n) \}$$
(1.20)

$$\sigma_{\bar{X}}^2 = \frac{1}{n^2} \{ \sigma_{\bar{X}}^2 + \sigma_{\bar{X}}^2 + \dots + \sigma_{\bar{X}}^2 \}$$
(1.21)

$$\sigma_{\bar{X}}^2 = \frac{1}{n^2} \{ n \sigma_X^2 \} = \frac{\sigma_X^2}{n}$$
(1.22)

Graphically, this result is shown in Figure 1.3. The distribution of X has a higher variance (it is more dispersed) than the distribution of  $\bar{X}$ .

## 1.7 Unbiasedness and efficiency

## 1.7.1 Unbiasedness

Because estimators are random variables, we can take expectations of the estimators. If the expectation of the estimator is equal to the true population parameter, then we say that this estimator is unbiased. Let  $\theta$  be the population parameter and let  $\hat{\theta}$  be a point estimator of  $\theta$ . Then,  $\hat{\theta}$  is unbiased if:

$$E(\hat{\theta}) = \theta \tag{1.23}$$

Example. The sample mean of *X* is an unbiased estimator of the population mean  $\mu_X$ :

#### 1 Random Variables, Sampling and Estimation



**Fig. 1.3** Probability density functions of *X* and  $\bar{X}$ .

$$E(\bar{X}) = E(\frac{1}{n}\sum_{i=1}^{n} x_i) = \frac{1}{n}E(\sum_{i=1}^{n} x_i)$$

$$= \frac{1}{n}\sum_{i=1}^{n}(E(x_i)) = \frac{1}{n}\sum_{i=1}^{n}\mu_X = \frac{1}{n}n\mu_X = \mu_X$$
(1.24)

**Unbiased estimator**. An estimator is unbiased if its expected value is equal to the true population parameter.

The bias of an estimator is just the difference between its expected value and the true population parameter:

$$Bias(\hat{\theta}) = E(\hat{\theta}) - \theta \tag{1.25}$$

## 1.7.2 Efficiency

It is not only important that an estimator is on average correct (unbiased), but also that it has a high probability of being close to the true parameter. When comparing two estimators,  $\hat{\theta}_1$  and  $\hat{\theta}_2$ , we say that  $\hat{\theta}_1$  is more efficient if  $var(\hat{\theta}_1) < var(\hat{\theta}_2)$ . A comparison of the efficiency between these two estimators in presented in Figure 1.4. The estimator with higher variance,  $(\hat{\theta}_2)$ , is more dispersed.



**Fig. 1.4** Efficiency of estimators  $\hat{\theta}_1$  and  $\hat{\theta}_2$ , with  $var(\hat{\theta}_1) < var(\hat{\theta}_2)$ .

**Most efficient estimator**. The estimator with the smallest variance from all unbiased estimators.

## 1.7.3 Unbiasedness versus efficiency

Both, unbiasedness and efficiency, are desired properties of an estimator. However, there may be conflicts in the selection between two estimators  $\hat{\theta}_1$  and  $\hat{\theta}_2$ , if, for example,  $\hat{\theta}_1$  is more efficients, but it is also biased. This case is presented in Figure 1.5.

The simplest way to select between these two estimators is to pick the one that yields the smallest mean square error (MSE):

$$MSE(\hat{\theta}) = var(\hat{\theta}) + bias(\hat{\theta})^2$$
(1.26)

#### **1.8** Estimators for the variance, covariance, and correlation

While we have already seen the populations formulas for the variance, covariance and correlation, it is important to keep in mind that we do not have the whole population. The data sets we will be working with are just samples of the populations. The formula for the sample variance is:

#### 1 Random Variables, Sampling and Estimation



**Fig. 1.5**  $\hat{\theta}_2$  is unbiased, but  $\hat{\theta}_1$  is more efficient.

$$s_X^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2$$
(1.27)

Notice how we changed the notation from  $\sigma^2$  to  $s^2$ . The first one denotes the population variance, while the second one refers to the sample variance. An estimator for the population covariance is given by:

$$s_{XY} = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{X}) (y_i - \bar{Y}).$$
(1.28)

Finally, the formula for the correlation coefficient,  $r_{XY}$ , is:

$$r_{XY} = \frac{\sum_{i=1}^{n} (x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{\sum_{i=1}^{n} (x_i - \bar{X})^2 \sum_{i=1}^{n} (y_i - \bar{Y})^2}}.$$
(1.29)

## **1.9** Asymptotic properties of estimators

Asymptotic properties of estimators just refers to their properties when the number of observations in the sample grows large and approached to infinity.



Fig. 1.6 The estimator is biased for small samples, but consistent.

## 1.9.1 Consistency

An estimator  $\hat{\theta}$  is said to be consistent if its bias becomes smaller as the sample size grows large. Consistency is important because many of the most common estimators used in econometrics are biased, then the minimum we should expect from these estimators is that the bias becomes small as we are able to obtain larger data sets. Figure 1.6 illustrates the concept of consistency by showing how an estimator of the population parameter  $\theta$  becomes unbiased as  $n \to \infty$ .

## 1.9.2 Central limit theorem

Having normally distributed random variables is important because we can then construct, for example, confidence intervals for its mean. However, what if a random variable does not follow a normal distribution? The central limit theorem gives us the answer.

**Central limit theorem**. States the conditions under which the mean of a sufficiently large number of independent random variables (with finite mean and variance) will be approximate a normal distribution.

Hence, even if we do not know the underlying distribution of a random variable, we will still be able to construct confidence intervals that will be approximately valid. In a numerical example, let's assume that the random variable X follows a

#### 1 Random Variables, Sampling and Estimation



Fig. 1.7 Distribution of the sample mean of a uniform distribution.

uniform distribution [-0.5,0.5]. Hence, it is equally likely that this random variable takes any value within this range. Figure 1.7 shows the distribution of the average of this random variable for n = 10, 20, and 100. All of these three distributions look very close to a normal distribution.

## Chapter 2 Simple Linear Regression

## 2.1 Simple linear model

The simple linear regression model shows how one known dependent variable is determined by a single explanatory variable (regressor). Is is written as:

$$Y_i = \beta_1 + \beta_2 X_i + u_i. \tag{2.1}$$

The subscript *i* refers to the observation i = 1, 2, ..., n, and  $Y_i$  is the dependent variable. We break down  $Y_i$  into two components, the deterministic (nonrandom) component  $\beta_1 + \beta_2 X_i$  and the stochastic (random) component  $u_i$ . The explanatory variable is  $X_i$  and the population parameters we want to estimate are given by intercept  $\beta_1$  and the slope  $\beta_2$ . The term  $u_i$  is the disturbance term. Figure 2.1 shows a graphical representation of the problem. The regression line  $Y_i = \beta_1 + \beta_2 X_i + u_i$  is shown as the upward sloping blue line. Only a single observation point at  $(X_i, Y_i)$  is illustrated. We can see how for this observation *i*, we break down  $Y_i$  into the disturbance term  $u_i$  given by the vertical distance between  $Y_i$  and  $\hat{Y}_i$  and the height of the regression line at point  $X_i$ , given by  $\beta_1 + \beta_2 X_i$ .

### 2.2 Least squares regression

The main idea in econometric analysis is to estimate the parameters  $\beta_1$  and  $\beta_2$ . The most popular estimator for these population parameters is the Ordinary Least Squares (OLS) estimator. Let the OLS estimators of  $\beta_1$  and  $\beta_2$  be  $b_1$  and  $b_2$ , respectively. Then, the fitter regression equation is:

$$Y_i = b_1 + b_2 X_i + e_i. (2.2)$$

The difference between Equations 2.1 and 2.2 is that the first correspond to the population, while the second is the sample counterpart. The idea in the OLS es-



**Fig. 2.1** Regression line  $Y_i = \beta_1 + \beta_2 X_i + u_i$ .

timator is simple, we want to pick values for the intercept  $b_1$  and slope  $b_2$  coefficients that are as close as possible to the actual data points. That is, we want to  $e_i$  ( $e_i = Y_1 - b_1 - b_2 X_i$ ) to be small. Because some of the  $e_i$  are positive and some are negative, we will first square them to have all positive numbers. Then, to take into account all data points we will sum across all observations. That is how our objective is to pick  $b_1$  and  $b_2$  to minimize the following residual sum of squares:

$$RSS = e_1^2 + e_2^2 + \dots + e_n^2 = \sum_{i=1}^n e_i^2.$$
 (2.3)

This minimization exercise yields the OLS estimators:

$$b_2 = \frac{\sum_{i=1}^{n} (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^{n} (X_i - \bar{X})^2}$$
(2.4)

for the slope coefficient, and

$$b_1 = \bar{Y} - b_2 \bar{X} \tag{2.5}$$

for the intercept. The derivation of the least squares coefficient estimators (Equations 2.4 and 2.5) has the following steps. We start with the regression equation:

$$Y_i = b_1 + b_2 X_i + e_i$$
$$\hat{Y}_i = b_1 + b_2 X_i$$

For observation i we obtain the residual, then square it and finally sum across all observations to obtain the residual sum of squares:

$$e_{i} = Y_{i} - \hat{Y}_{i}$$

$$e_{i}^{2} = (Y_{i} - \hat{Y}_{i})^{2}$$

$$\sum_{i=1}^{n} e_{i}^{2} = \sum_{i=1}^{n} (Y_{i} - \hat{Y}_{i})^{2}$$
(2.6)

The coefficients b<sub>1</sub> and b<sub>2</sub> are chosen to minimize the residuals sum of squares:

$$\min_{b_1, b_2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$\min_{b_1, b_2} \sum_{i=1}^n (Y_i - b_1 - b_2 X_i)^2$$
(2.7)

The first order necessary condition are:

$$-2\sum_{i=1}^{n} (Y_i - b_1 - b_2 X_i) = 0 \quad \text{w.r.t.} \quad b_1$$
 (2.8)

$$-2\sum_{i=1}^{n} X_i(Y_i - b_1 - b_2 X_i) = 0 \quad \text{w.r.t.} \quad b_2$$
 (2.9)

Dividing Equation 2.9 by n and working through some math we obtain the OLS estimators for the constant:

$$b_1 = \bar{Y} - b_2 \bar{X}.$$

Plugging this result into Equation 2.9 we obtain:

$$b_2 = \frac{\sum_{i=0}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=0}^n (X_i - \bar{X})^2}.$$

## 2.3 Interpretation of the regression coefficients

If the estimated regression equation is given by:

$$\widehat{wage}_i = 4.64 + 0.09 exper_i,$$
 (2.10)

where *wage* is the hourly wage measured in dollars, and *exper* is the number of years of experience, then the interpretation of the slope coefficient is the following:

$$\frac{\Delta wage}{\Delta exper} = 0.09.$$

Therefore, if the change in the number of years of experience is one,  $\Delta exper$ , then the change in the hourly wage in dollars is given by  $\Delta wage = 0.09$ . In words, an additional year of experience will increase your hourly wage by 0.09 dollars (or 9 cents). For the interpretation of the intercept, just consider the case where someone has not experience, exper = 0. Then, this person's predicted wage will be 4.64 dollars.

If the estimated regression equation takes the form:

$$\log wage_i = 1.38 + 0.02 exper_i,$$
 (2.11)

where the log *wage* is the natural logarithm of *wage*, then the interpretation is different. Here, if the number of years of experience increases by one, the wage increases by  $2\% (0.02 \times 100 \text{ percent})$ . Finally, for the following estimated equation:

$$\log wage_i = 0.98 + 0.26 \log exper_i.$$
(2.12)

A one percent increase in *exper* will increase *wage* by 0.25 percent. The 0.26 is interpreted as an elasticity.

## 2.4 Goodness of fit

How good is the regression equation in explaining the variation in variable Y? First we need a way to measure the total variation in Y. Let's try the sum of squared deviations about the sample mean of Y. That is,

$$\sum_{i=1}^{n} (Y_i - \bar{Y})^2 \tag{2.13}$$

Now, let's start with a simple equality:

$$Y_i - \bar{Y} = Y_i - \bar{Y}.$$

If we add and subtract  $\hat{Y}_i$  on the right hand side of the above equality, we have

$$Y_i - \bar{Y} = Y_i - \bar{Y} + \hat{Y}_i - \hat{Y}_i$$
  
$$Y_i - \bar{Y} = (\hat{Y}_i - \bar{Y}) + (Y_i - \hat{Y}_i)$$

Squaring both sides of the equation and then summing across all observations *i* we obtain:

$$\sum_{i=1}^{n} (Y_i - \bar{Y})^2 = \sum_{i=1}^{n} (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$$
(2.14)

$$TSS = ESS + RSS. \tag{2.15}$$



**Fig. 2.2** Decomposition of  $\hat{Y}_i - \bar{Y}$ .

Notice that the sum of deviations from the mean is zero, that is why there are only two components on the right hand side. The *TSS* is the Total Sum of Squares, as presented in Equation 2.13. The first term on the right hand side is ESS, the Explained Sum of Squares, and the second term on the right hand side is the RRS, Residual Sum of Squares. This decomposition of the variable *Y* into two components can be appreciated in Figure 2.2. For every observation  $Y_i$  in the sample, the distance between  $Y_i$  and  $\bar{Y}$  can be decomposed in two, the part that the regression equation can explain,  $\hat{Y}_i - \bar{Y}$ , and the part that the regression equation cannot explain,  $Y_i - \hat{Y}_i$ .

What is the proportion of the variation in *Y* that is explain by the regression equation? We just need to divide Equation 2.15 by *TSS* and define the ratio of *ESS* to *TSS* as the proportion of the explained variation in *Y*, the  $R^2$ :

$$1 = \frac{ESS}{TSS} + \frac{RSS}{TSS}$$
(2.16)

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$$
(2.17)

$$R^{2} = \frac{\sum_{i=1}^{n} (\hat{Y}_{i} - \bar{Y})^{2}}{\sum_{i=1}^{n} (Y_{i} - \bar{Y})^{2}} = 1 - \frac{\sum_{i=1}^{n} e_{i}^{2}}{\sum_{i=1}^{n} (Y_{i} - \bar{Y})^{2}}$$
(2.18)

The  $R^2$  is a number between zero and one, being higher when the model explains more of the variation in Y. Figures 2.3 and 2.4 illustrate how the regression line explain the variation in Y when the  $R^2$  is low and high, respectively.



**Fig. 2.3** Low  $R^2$ .



**Fig. 2.4** High  $R^2$ .

## Chapter 3 Properties and Hypothesis Testing

## 3.1 Types of data

The regression techniques developed in previous chapters can be applied to three different kinds of data.

- 1. Cross-sectional data.
- 2. Time series data.
- 3. Panel data.

The first consists on observing various economic unit (e.g. firms, countries, households, individuals) at one point in time. For example, we observe the wages, experience and education of many individuals, only once and at all at the same time. The second consists on observing the same economic unit at different point in time. For example, we observe daily stock prices over many years. Finally, the third combines the characteristics of the first and the second. That is, we observe various economic units at repeated points in time. For example, we have information about the inflation, unemployment and GDP of a group of countries and over many years.

## 3.2 Assumptions of the model

When the regressors in our econometric model are non stochastic, we will make the following six assumptions.

1. The model is linear in the parameters and it is correctly specified.

$$Y = \beta_1 + \beta_2 X + u \tag{3.1}$$

$$Y = \beta_1 X^{\beta_2} + u \tag{3.2}$$

Equation 2.1 is linear in  $\beta$ , while Equation 2.2 is not.

2. There is some variation in the regressor in the sample. We need variation in the variable *X* to identify the relationship. Consider the OLS estimator for  $\beta_2$ :

$$b_2 = \frac{\sum_{i=0}^{n} (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=0}^{n} (X_i - \bar{X})^2}.$$
(3.3)

If there is no variation in X, then the denominator is zero and we cannot obtain  $b_2$ .

3. The expected value of the disturbance term is zero.

$$E(u_i) = 0 \quad \text{for all} \quad i. \tag{3.4}$$

Some  $u_i$  will be negative, some will be positive, but on average they will be zero. If a constant is included in the model, the condition is satisfied automatically.

4. The disturbance term is homoscedastic.

*Homoscedasticity* means that the variance of the error terms  $u_i$  is constant across all observations *i*. Hence, we can write:

$$\sigma_{u_i}^2 = \sigma_u^2 \quad \text{for all} \quad i. \tag{3.5}$$

Because the error term has zero mean (from assumption 3), then the population variance of  $u_i$  is equal to:

$$E(u_i^2) = \sigma_u^2 \quad \text{for all} \quad i. \tag{3.6}$$

 $\sigma_u^2$  is a population parameter, therefore it is unknown and need to be estimated. 5. The values of the disturbance terms have independent distributions.

$$u_i$$
 is distributed independently of  $u_j$  for all  $j \neq i$ . (3.7)

This means that there is no *autocorrrelation* in the error term. This means that the population covariance between  $u_i$  and  $u_j$  is zero:

$$\sigma_{u_i u_i} = 0. \tag{3.8}$$

With assumptions 1 through 5, we says that OLS coefficients are BLUE: Best Linear Unbiased Estimators. They are best, because they have the smallest variance across all unbiased estimators.

6. The disturbance term has a normal distribution.

$$u_i \sim N[0, \sigma_u^2]$$
 for all *i*. (3.9)

The error term is distributed normal with mean zero and variance  $\sigma_u^2$ . This assumption becomes useful at the time of performing *t* tests, *F* tests, and constructing confidence intervals for  $\beta_1$  and  $\beta_2$  using the regression results. The justification for this assumption depends on the *central limit theorem*. This one state that if a random variable is the composite result of the effects of a large number of

24

3.4 Precision of the coefficients

other random variables (that are not necessarily normal), it will have an approximately normal distribution.

## 3.3 Unbiasedness of the coefficients

Recall that an estimator  $\hat{\theta}$  is unbiased if  $E(\hat{\theta}) = \theta$ . The expected value of the estimator is equal to the true population parameter. For the slope coefficient in the OLS regression we have:

$$b_{2} = \frac{\sum_{i=0}^{n} (X_{i} - \bar{X})(Y_{i} - \bar{Y})}{\sum_{i=0}^{n} (X_{i} - \bar{X})^{2}}$$

$$= \beta_{2} + \frac{\sum_{i=0}^{n} (X_{i} - \bar{X})u_{i}}{\sum_{i=0}^{n} (X_{i} - \bar{X})^{2}}$$

$$= \beta_{2} + \sum_{i=1}^{n} a_{i}u_{i}$$
(3.10)

where

$$a_i = \frac{(X_i - \bar{X})}{\sum_{i=0}^n (X_i - \bar{X})^2}.$$
(3.11)

Thus, this shows that  $b_2$  is equal to its true value,  $\beta_2$ , plus a linear combination of the values of the error terms. If we take expectations of  $b_2$  we have:

$$E(b_2) = E(\beta_2) + E\left(\sum_{i=1}^n a_i u_i\right) = \beta_2 + \sum_{i=1}^n E(a_i u_i) = \beta_2 + \sum_{i=1}^n a_i E(u_i) = \beta. \quad (3.12)$$

The term  $a_i$  goes out of the expectation because  $a_i$  is only a function of constant *X*s. In addition, the last equality holds because  $E(u_i) = 0$ . Hence,  $b_2$  is an unbiased estimator of  $\beta_2$ ,  $E(b_2) = \beta_2$ .

## 3.4 Precision of the coefficients

We are also interested on how precise  $b_1$  and  $b_2$  are in estimating the population parameters  $\beta_1$  and  $\beta_2$ . A measure of this precision are their population variances, given by:

$$\sigma_{b_1}^2 = \sigma_u^2 \left( \frac{1}{n} + \frac{\bar{X}}{\sum_{i=0}^n (X_i - \bar{X})^2} \right), \text{ and }$$
(3.13)

$$\sigma_{b_2}^2 = \frac{\sigma_u^2}{\sum_{i=0}^n (X_i - \bar{X})^2}$$
(3.14)

One concern in the implementation of the above formulas is that  $\sigma_u^2$  is an unknown population parameter and need to be estimated. A natural estimator for this regression variance is the variance of the regression errors. Because the population regression errors  $u_i$  are also unknown, we use the sample counterparts  $e_i$  and adjust for the corresponding degrees of freedom. Hence, we have:

$$S_u^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2.$$
(3.15)

This  $S_u^2$  is the unbiased estimator of  $\sigma_u^2$ , and n-2 are the degrees of freedom. We subtract two from the sample size because we are estimating two parameters: the regression constant and one slope coefficient. Then, we use the following formulas to estimate the standard errors of  $b_1$  and  $b_2$ :

$$S_{b_1} = \sqrt{S_u^2 \left(\frac{1}{n} + \frac{\bar{X}}{\sum_{i=0}^n (X_i - \bar{X})^2}\right)}, \text{ and } (3.16)$$

$$S_{b_2} = \sqrt{\frac{S_u^2}{\sum_{i=0}^n (X_i - \bar{X})^2}}.$$
(3.17)

## 3.5 The Gauss-Markov theorem

The *Gauss-Markov theorem* simply states that when assumptions 1 through 5 above are satisfied, the OLS estimators are Best Linear Unbiased Estimators (BLUE) of the regression parameters. Best refers to smallest variance.

## 3.6 Hypotheses testing

*Hypothesis testing* is simply a method of making decisions using data. It starts with the formulation of the null and the alternative hypotheses and then uses some test statistics to assess the truth of the null hypothesis.

## 3.6.1 Formulation of the null hypothesis

The formulation of the null hypothesis starts with a relationship in mind. For example, that the percentage rate of price inflation (p) depends on the percentage rate of wage inflation (w) following the linear equation:

$$p_i = \beta_1 + \beta_2 w_i + u_i \tag{3.18}$$

#### 3.6 Hypotheses testing

Then, you want to test the hypothesis that the price inflation is equal to the wage inflation. This is denoted by  $H_0$  and it is know as the *null hypothesis*. In addition, we also define an alternative hypothesis, denoted by  $H_1$  and represents the conclusion of the test if the null hypothesis is rejected. For our example the null and the alternative hypothesis are written as:

$$H_0: \beta_2 = 1 \tag{3.19}$$

$$H_1: \beta_2 \neq 1 \tag{3.20}$$

In general, the null and alternative hypotheses are:

$$H_0: \beta_2 = \beta_2^0 \tag{3.21}$$

$$H_1: \beta_2 \neq \beta_2^0.$$
 (3.22)

### 3.6.2 t-tests

Recall that  $\beta_2$  is unknown and that we have to use the estimate  $b_2$ . Then, the decision rule to reject the null hypothesis should compare the estimate  $b_2$  with the hypothesized value  $\beta_2^0$ . Intuitively, if the values are far apart, then there is evidence against the null. This comparison should take into account the fact that  $b_2$  is subject to some sampling variation (it is not the actual  $\beta_2$ ). We will use the following statistic:

$$z = \frac{b_2 - \beta_2^0}{\sigma_{b_2}}$$
(3.23)

The numerator is just the distance between the regression estimate and the hypothesized value, with the denominator is the standard deviation of  $b_2$ , given by the square root of the expression in Equation 3.14. z is the number of standard deviations between  $b_2$  and  $\beta_2$ . For a known  $\sigma_{b_2}$ , this one follows a normal distribution. However  $\sigma_{b_2}$  is unknown and we need to use the estimate of the standard error of  $b_2$ . This one is given by  $S_{b_2}$  and it is presented in Equation 3.17. Then we use the following *t*-statistic:

$$t = \frac{b_2 - \beta_2^0}{S_{b_2}} \tag{3.24}$$

To know if the deviations between  $b_2$  and  $\beta_2^0$  are significantly large, we compare this *t*-statistic with the critical values from the table *t* distribution with n - 2 degrees of freedom. The null hypothesis is not rejected if the following condition is met:

$$-t_{n-2,\alpha/2} \le \frac{b_2 - \beta_2^0}{S_{b_2}} \le t_{n-2,\alpha/2}$$
(3.25)

Where  $t_{n-2,\alpha/2}$  is just the notation of the critical value than comes from the *t* distribution with n-2 degrees of freedom and at significance level  $\alpha$ . The *significance* 



Fig. 3.1 Acceptance region for the *t*-test.

*level* is the probability that we reject the null hypothesis when in fact it is true. The rejection regions are illustrated in Figure 3.1.

## 3.6.3 Confidence intervals

The *confidence interval* indicates the reliability of an estimate. The confidence interval for the population parameter  $\beta_2$  can be derived from Equation 3.25 in the following way:

$$1 - \alpha = P\left(-t_{n-2,\alpha/2} \le \frac{b_2 - \beta_2}{S_{b_2}} \le t_{n-2,\alpha/2}\right)$$
(3.26)  

$$1 - \alpha = P\left(-t_{n-2,\alpha/2} \cdot S_{b_2} \le b_2 - \beta_2 \le t_{n-2,\alpha/2} \cdot S_{b_2}\right)$$
  

$$1 - \alpha = P\left(b_2 - t_{n-2,\alpha/2} \cdot S_{b_2} \le \beta_2 \le b_2 + t_{n-2,\alpha/2} \cdot S_{b_2}\right)$$

The meaning of the above equation is that the population parameter  $\beta_2$  will be between the lower confidence limit  $b_2 - t_{n-2,\alpha/2} \cdot S_{b_2}$  and the upper confidence limit  $b_2 + t_{n-2,\alpha/2} \cdot S_{b_2}$  with probability  $(1 - \alpha)$  or  $100 \times (1 - \alpha)\%$ . The *p* values provide an alternative approach to reporting the significance of regression coefficients or when carrying out more general hypothesis testing. As you can see from Equation 3.25 and Figure 3.1, different significance levels  $\alpha$  can yield a different conclusion in the rejection or not of the null hypothesis. The p value of a hypothesis test represent the minimum significance level at which the null is rejected. Then, when the p value is below the significance level  $\alpha$  we reject the null.



**Fig. 3.2** Confidence interval for  $\beta_2$ .

## 3.6.4 F test

A useful tool if we want to test if there is no relationship between X and Y if the F test. In the simple linear regression model with only one slope coefficient, the null and the alternative in an F test are:

$$H_0: \beta_2 = 0 \tag{3.27}$$

$$H_1: \beta_2 \neq 0.$$
 (3.28)

This test is build on the idea of testing how good is the regression model in explaining the variation in Y. In Equation 2.15 we already separated the variation of Y into its 'explained' and 'unexplained' components. These are:

$$\sum_{i=1}^{n} (Y_i - \bar{Y})^2 = \sum_{i=1}^{n} (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$$
(3.29)

$$TSS = ESS + RSS. \tag{3.30}$$

The total sum of squares (TSS) is the summation of the explained sum of squares (ESS) and the residual sum of squares (RSS). Then, the *F statistic* for goodness of fit of a regression is written as the explained sum of squares, per explanatory variable, divided by the residual sum of squares, per remaining degrees of freedom:

$$F = \frac{ESS/(k-1)}{RSS/(n-k)}$$
(3.31)

SUMMARY OUTPUT

Regression S	tatistics					
Multiple R	0.2346					
R Square	0.0551					
Adjusted R Square	0.0543					
Standard Error	4.5323					
Observations	1260					
ANOVA						
	df	SS	MS	F	Significance F	
Regression	1	1505.5387	1505.5387	73.2906	0.0000	
Residual	1258	25841.9006	20.5421			
Total	1259	27347.4393				
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%
Intercept	4.6425	0.2326	19.9615	0.0000	4.1862	5.0988
X Variable 1	0.0914	0.0107	8.5610	0.0000	0.0705	0.1124

Fig. 3.3 Regression output in MS Excel.

where k is the total number of coefficients we are estimating, hence (k - 1) is the number of slope coefficients. That is, the total number of parameters we are estimating minus the constant parameter. If we divide the numerator and the denominator by *TSS*, then the *F* statistics can be written in terms of the  $R^2$  as follows:

$$F = \frac{(ESS/TSS)/(k-1)}{(RSS/TSS)/(n-k)} = \frac{R^2/(k-1)}{(1-R^2)/(n-k)}$$
(3.32)

If this *F* statistic is greater that the critical value from the table *F* distribution with (k-1) and (n-k) degrees of freedom,  $F_{k-1,n-k}$ , we reject the null hypothesis and conclude that the regression model does not significantly explain the variation in variable *Y*. For the simple regression model with only one slope coefficient, k = 2, we have:

$$F = \frac{R^2}{(1 - R^2)/(n - 2)}.$$
(3.33)

If this *F* statistic >  $F_{1,n-2}$  we reject the null hypothesis presented in Equation 3.28.
#### 3.7 Computer output

# 3.7 Computer output

The computer regression output is very similar across different statistical packages. Figure 3.3 shows the output using MS Excel for the estimation of the following simple regression model:

$$wage_{=}\beta_{1} + \beta_{2}exper_{i} + u_{i} \tag{3.34}$$

To obtain the regression estimated coefficients we use Equations 2.4 and 2.5:

$$b_2 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = 0.091$$
(3.35)

$$b_1 = \bar{Y} - b_2 \bar{X} = 4.642 \tag{3.36}$$

The total sum of squares, estimates sum of squares, and residual sum of squares are obtained using 2.15 and 2.15:

$$TSS = \sum_{i=1}^{n} (Y_i - \bar{Y})^2 = 27347.439$$
(3.37)

$$ESS = \sum_{i=1}^{n} (\hat{Y}_i - \bar{Y})^2 = 1505.539$$
(3.38)

$$RSS = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 = 25841.901$$
(3.39)

The regression  $R^2$  comes from Equation 2.18:

$$R^{2} = 1 - \frac{\sum_{i=1}^{n} e_{i}^{2}}{\sum_{i=1}^{n} (Y_{i} - \bar{Y})^{2}} = 0.055$$
(3.40)

From the square root of Equation 3.15:

$$S_u = \sqrt{\frac{1}{n-2} \sum_{i=1}^n e_i^2} = 4.532 \tag{3.41}$$

Then, the standard errors of the coefficients are computer using Equations 3.17 and 3.17:

$$S_{b_1} = \sqrt{S_u^2 \left(\frac{1}{n} + \frac{\bar{X}}{\sum_{i=0}^n (X_i - \bar{X})^2}\right)} = 0.233$$
(3.42)

$$S_{b_2} = \sqrt{\frac{S_u^2}{\sum_{i=0}^n (X_i - \bar{X})^2}} = 0.011$$
(3.43)

The *F* statistic uses Equation 3.32:

3 Properties and Hypothesis Testing

$$F = \frac{R^2/(k-1)}{(1-R^2)/(n-k)} = 73.291$$
(3.44)

The *t* statistics use Equation 3.24:

$$t = \frac{b_1}{S_{b_1}} = 19.961 \tag{3.45}$$

$$t = \frac{b_2}{S_{b_2}} = 8.561 \tag{3.46}$$

Finally, for the 95% upper and lower confidence levels, we use Equation 3.26:

$$b_1 - t_{n-2,\alpha/2} \cdot S_{b_1} = 4.186 \tag{3.47}$$

$$b_1 + t_{n-2,\alpha/2} \cdot S_{b_1} = 5.099 \tag{3.48}$$

$$b_2 - t_{n-2,\alpha/2} \cdot S_{b_2} = 0.071 \tag{3.49}$$

$$b_2 + t_{n-2,\alpha/2} \cdot S_{b_2} = 0.112 \tag{3.50}$$

# Chapter 4 Multiple Regression Analysis

The simple linear regression covered in Chapter 2 can be generalized to include more than one variable. Multiple regression analysis is an extension of the simple regression analysis to cover cases in which the dependent variable is hypothesized to depend on more than one explanatory variable. While much of the analysis is an extension of the simple case, we have two main complications. (1) We need to discriminate between the effects of one variable and the effects of the other explanatory variables. (2) We have to decide which variables to include in the regression equation. In this chapter we will focus on the extension of the linear regression model and in (1). In a later chapter we will discuss (2).

# 4.1 Interpretation of the coefficients

Consider the following population multiple regression model with (k-1) regressors:

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k + u.$$
(4.1)

A simple example of a multiple regression model is:

$$CRIME_i = \beta_1 + \beta_2 POPULATION_i + \beta_3 UNEMPLOY_i + \beta_4 POLICE_i + u_i, \quad (4.2)$$

where *i* refers to the city, *CRIME* is crime rates, *POPULATION* is just the number people in city *i*, *UNEMPLOY* is the unemployment rate, and *POLICE* is the number of police officers. To estimate the  $\beta$ s in Equation 4.2 you may need to observe crime rates and all the other variables for *n* cities. As before, *u* is the disturbance term. Because we have more that one regressor, the simple two dimensional characterization illustrated in Figure 2.1 is no longer applicable. Now, we have a (k - 1) dimensional problem. In our crime example we would need to have a 4D graph!

The sample counterpart of Equation 4.2 is:

$$CRIME_i = b_1 + b_2 POPULATION_i + b_3 UNEMPLOY_i + b_4 POLICE_i + e_i$$
, (4.3)

where the *bs* are the sample estimates of the  $\beta$ s, and are estimated using computer software via Ordinary Least Squares. We also express this relationship as the 4D fitted plane:

$$CRIME_i = b_1 + b_2 POPULATION_i + b_3 UNEMPLOY_i + b_4 POLICE_i.$$
(4.4)

Notice that we no longer write the disturbance term. Moreover,  $CRIME_i$  is the fitted or predicted value of  $CRIME_i$ . The interpretation of the coefficients is the same as before. If the number of police officers increases by one, then the crime rate will change by  $b_4$ . Similar interpretation follows for  $b_2$  and  $b_3$ .

### 4.2 Ordinary Least Squares

The OLS estimates are obtained in the same fashion as before. The unknown relationship is given by:

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k X_{ki} + u_i.$$
(4.5)

The fitted OLS regression is:

$$\hat{Y}_i = b_1 + b_2 X_{2i} + b_3 X_{3i} + \dots + b_k X_{ki}.$$
(4.6)

Then, the OLS regression residuals are:

$$e_i = Y_i - \hat{Y} = Y_i - b_1 - b_2 X_{2i} - b_3 X_{3i} - \dots - b_k X_{ki}.$$
(4.7)

Recall that OLS minimizes the sum of squared residuals

$$\min_{b_1, b_2, \dots, b_k} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2, \tag{4.8}$$

where  $RSS = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$  is the sum of squared residuals. We need to take the derivative of the *RSS* with respect to  $b_1, b_2, ..., b_k$  and obtain *k* first order conditions. This yields a system of *k* equations with *k* unknowns, where the solution is the OLS estimators of the  $\beta$ s.

# 4.3 Assumptions

1. The model is linear in the parameters and correctly specified

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k + u.$$
(4.9)

- 4.4 Properties of the coefficients
- 2. There is no exact linear relationship among the regressors in the sample. This is called multicollinearity.
- 3. The disturbance term has expectation zero

$$E(u_i) = 0 \qquad \text{for all} \quad i. \tag{4.10}$$

4. The disturbance term is homoscedastic.

$$\sigma_{u_i}^2 = \sigma_u^2 \qquad \text{for all} \quad i. \tag{4.11}$$

5. The values of the disturbance term have independent distributions.

$$u_i$$
 is distributed independently of  $u_{i'}$  for all  $i' \neq i$ . (4.12)

6. The distribution term has a normal distribution.

$$u_i \sim N[0, \sigma^2]$$
 for all *i*. (4.13)

All the Xs are nonstochastic.

## **4.4 Properties of the coefficients**

### 4.4.1 Unbiasedness

The OLS estimator  $b_j$  of  $\beta_j$  is unbiased:

$$E(b_j) = \beta_j \tag{4.14}$$

# 4.4.2 Efficiency

Following the results from the Gauss-Markov theorem, we have that OLS yields the most efficient linear estimators, in the sence that they are the one with the smallest variance among all linear estimators.

# 4.4.3 Precision of the coefficient, t tests, and confidence intervals

Beside our interest on the point estimates, we are also interested in performing hypotheses testing and building confidence intervals. To do this we need a measure of the precision of the coefficients. While we will not show the derivation here (as it required matrix algebra), each of the  $b_j$  has an standard error,  $S_{b_j}$ .

The null and alternative hypotheses about population coefficient *j* is written as:

$$H_0: \beta_j = \beta_j^0 \tag{4.15}$$

$$H_1: \beta_j \neq \beta_j^0. \tag{4.16}$$

which can be tested using the following *t*-statistic:

$$t = \frac{b_j - \beta_j^0}{S_{b_j}}$$
(4.17)

The null is not rejected if the following condition is met:

$$-t_{n-k,\alpha/2} \le \frac{b_j - \beta_j^0}{S_{b_j}} \le t_{n-k,\alpha/2}$$
(4.18)

Notice the difference between Equation 4.18 and Equation 3.25. The critical value from the *t* distribution,  $t_{n-k,\alpha/2}$ , now has n-k degrees of freedom because we are estimating *k* parameters, rather than just 2 as in the simple regression model. The intuition behind Figures 3.1 and 3.1 still hold. The computer software will also give you the p-value associated with the *t* test. If the p-value is below your  $\alpha$ , you reject the null hypothesis.

For the construction of the confidence intervals we have:

$$1 - \alpha = P\left(-t_{n-k,\alpha/2} \le \frac{b_j - \beta_j}{S_{b_j}} \le t_{n-k,\alpha/2}\right)$$

$$1 - \alpha = P\left(-t_{n-k,\alpha/2} \cdot S_{b_j} \le b_j - \beta_j \le t_{n-k,\alpha/2} \cdot S_{b_j}\right)$$

$$1 - \alpha = P\left(b_j - t_{n-k,\alpha/2} \cdot S_{b_j} \le \beta_j \le b_j + t_{n-k,\alpha/2} \cdot S_{b_j}\right).$$
(4.19)

## 4.5 Regression output in Gretl

Gretl is an open-source (free) software package for econometric analysis written in the C programming language. It can be downloaded from:

http://gretl.sourceforge.net/

Just follow the instructions to install it in your computer.

Once you loaded the data set in Gretl, to estimate Equation 4.2 you need to go to  $Model \rightarrow Ordinary$  Least Squares. The regression output is:

#### 4.6 Multicollinearity

Model 1: OLS, using observations 1-92 Dependent variable: crimes

	coeffi	cient	sto	. er	ror	t-ratio	p-value	
const pop	2193.3	4 652716	3918 (	.06	6262	0.5598 6.143	0.5770 2.30e-08	***
unem officers	-279.2 15.0	91 406	407	.791	60	-0.6849 4.205	0.4952 6.25e-05	***
Mean depende Sum squared R-squared F(3, 88) Log-likeliho Schwarz crit	ent var resid pod cerion	39663. 1.39e+ 0.8272 140.51 -996.73 2011.5	53 10 93 07 10 49	S.D. S.E. Adju P-va Akai Hann	depe of r sted 1 lue(F ke cr an-Qu	ndent var egression R-squared ) iterion inn	29692.10 12548.04 0.821405 1.90e-33 2001.462 2005.533	

Excluding the constant, p-value was highest for variable 3 (unem)

A standard way to present the regression output is:

crimes = 
$$2193.34_{(3918.1)} - 279.291_{(407.79)}$$
 unem + 15.0406 officers + 0.0652716 pop  
(3.5766) (0.010626)  
 $N = 92$   $\bar{R}^2 = 0.8214$   $F(3,88) = 140.51$   $\hat{\sigma} = 12548.$   
(standard errors in parentheses)

To obtain the confidence intervals for the coefficients as presented in Equation 4.19 in the Gretl regression output window you need to go to Analysis  $\rightarrow$  Confidence intervals for the coefficients to obtain:

t(88, 0.025) = 1.987VARIABLE COEFFICIENT 95% CONFIDENCE INTERVAL 2193.34 -5592.97 9979.66 const 0.0652716 0.0441542 0.0863890 pop -279.291 -1089.69 531.107 unem officers 15.0406 7.93282 22.1483

# 4.6 Multicollinearity

*Multicollinearity* is when two explanatory variables are highly correlated. In addition, if their coefficients have a large population variances, we are at risk of getting erratic estimates of the coefficients. There could also be multicollinearity when there is an approximate linear relationship between more than two variables.

A simple test for multicollinearity is based in the Variance Inflation Factors. To implement this text in Gretl, in the regression output window go to Test  $\rightarrow$  Collinearity:

```
Variance Inflation Factors
```

Based on these results, we do not have a multicollinearity problem in the estimation of Equation 4.2.

# **4.7 Goodness of fit:** $R^2$ and $\bar{R}^2$

The  $R^2$  in multiple regression analysis has the same interpretation as in a simple regression. It is the proportion of the variation in *Y* explained by the regression model

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$$
(4.20)

$$R^{2} = \frac{\sum_{i=1}^{n} (\hat{Y}_{i} - \bar{Y})^{2}}{\sum_{i=1}^{n} (Y_{i} - \bar{Y})^{2}} = 1 - \frac{\sum_{i=1}^{n} e_{i}^{2}}{\sum_{i=1}^{n} (Y_{i} - \bar{Y})^{2}}$$
(4.21)

where  $\hat{Y}$  represents the fitted values of the regression equation

$$\hat{Y} = b_1 + b_2 X_2 + b_3 X_3 + \dots + b_k X_k. \tag{4.22}$$

## **4.8** *F* tests

Given the population regression model

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k + u,$$
(4.23)

we can use the *F* test to test if all the slope coefficients  $\beta_2, \beta_3, \ldots, \beta_k$  are jointly equal to zero. That is, let the null hypothesis be:

$$H_0: \beta_2 = \beta_3 = \dots = \beta_k = 0. \tag{4.24}$$

#### 4.9 Adjusted $R^2$ , $\bar{R}^2$

The alternative hypothesis  $(H_0)$  is that at least one of the slope coefficients is different from zero. The multiple regression version of the *F* statistic is:

$$F_{k-1,n-k} = \frac{ESS/(k-1)}{RSS/(n-k)}.$$
(4.25)

The idea is to compare this *F* statistic to the critical level found in the *F* distribution tables with k - 1 and n - k degrees of freedom. Computer software automatically computes this *F* statistic and the corresponding p-value for the null in Equation 4.22. This *F* statistic can also be written in terms of the  $R^2$ :

$$F_{k-1,n-k} = \frac{R^2/(k-1)}{(1-R^2)/(n-k)}.$$
(4.26)

Consider the example presented in Section 4.5. The *F* statistic is 140.5107 with 3 and 88 degrees of freedom and has a corresponding p-value of 0.000. Then, because the p-value is below  $\alpha = 5\%$  then we reject the null hypothesis that the slope coefficients on pop, unem, and officers are jointly equal to zero.

# **4.9** Adjusted $R^2$ , $\bar{R}^2$

One concern with the  $R^2$  is that it will always go up as we include more variables into the model. Hence, it is a poor way to compare models. On the other hand a similar statistic, the adjusted  $R^2$  ( $\overline{R}^2$ ) is built on the  $R^2$  but with the difference that  $\overline{R}^2$  penalizes for the loss of the degrees of freedom as we include more variables into the model. Therefore, the  $\overline{R}^2$  can either go up or down as we include more variable into the model. It is defined as:

$$\bar{R}^2 = R^2 - \frac{k-1}{n-k}(1-R^2).$$
(4.27)

# Chapter 5 Transformations of Variables and Interactions

## 5.1 Basic idea

One limitation in the linear regression analysis is that the dependent variable has to be linear in the parameters:

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k + u.$$
(5.1)

However, there are equations that are not linear, for example:

$$Y = \beta_1 + \beta_2 X_2 + X_3^{\beta_3} + u.$$
 (5.2)

This Equation 5.2 cannot be estimated using OLS. One way to estimate nonlinear models is by using Nonlinear Least Squares (NLS), which is an extension of the methods we discussed before. In this chapter, rather that focusing on NLS, we will see how transformations in the variables can allow us to use OLS on a variety of non-linear models. For example, consider the estimation of the following Cobb-Douglas production function:

$$P_i = A L_i^{\beta_2} K_i^{\beta_3} e^{\varepsilon_i}, \tag{5.3}$$

where  $P_i$  is total production or total output, A is a technology constant,  $K_i$  is the amount of capital, and  $L_i$  is labor. Taking natural logs we have:

$$\log P_i = \log A + \beta_2 \log L_i + \beta_3 \log K_i + \varepsilon_i.$$
(5.4)

If we simple set  $Y_i = \log P_i$ ,  $\beta_1 = \log A$ ,  $X_2 = \log L_i$ , and  $X_3 = \log K_i$  we can write Equation 5.4 as:

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \varepsilon_i, \qquad (5.5)$$

that can be easily estimated via OLS.  $\beta_2$  and  $\beta_3$  will correspond to the ones given in Equation 5.3. Another example of a model that can be estimated with OLS is:

$$Y_i = \beta_1 + \beta_2 Z_{2i}^2 + \beta_3 \sqrt{Z_{3i}} + \beta_4 \frac{1}{Z_{4i}} + \varepsilon_i.$$
(5.6)

5 Transformations of Variables and Interactions

We just need to replace  $X_{2i} = Z_{2i}^2, X_{3i} = \sqrt{Z_{3i}}, X_{4i} = \frac{1}{Z_{4i}}$ .

# 5.2 Logarithmic transformations

To explain the logarithm transformation let's go over one example in Gretl. If we want to estimate the following model:

 $\log \operatorname{crime}_{i} = \beta_{1} + \beta_{2} \log \operatorname{pop}_{i} + \beta_{3} \operatorname{unem}_{i} + \beta_{4} \operatorname{offi}_{i} + u_{i}, \qquad (5.7)$ 

you need to create the new variables first. Go to  $Add \rightarrow Define$  new variable and type:

logcrime = log(crime)

This will generate the new variable logorime. Do the same thing for log population and then estimate the model. The regression output is:

```
Model 1: OLS, using observations 1-92
        Dependent variable: logcrime
                                                        coefficient std. error t-ratio p-value
                 _____
                const-0.7097350.807193-0.87930.3817unem-0.004568480.00903041-0.50590.6142
              offi0.0001449150.00503041-0.50590.6142logpop0.8640440.066278212.04
                                                                                                                                                                                                          2.92e-022 ***
        Mean dependent var 10.33774 S.D. dependent var 0.742056
        Sum squared resid 6.883563
R-squared 0.862628
                                                                                                                              S.E. of regression 0.279683
Adjusted R-squared 0.857945
        Sum squared0.862628AugustR-squared0.862628AugustT2 081184.1989P-value(F)Abaike crit
       F(3, 88) 184.1989 F-Value (1,
Log-likelihood -11.28034 Akaike criterion 30.56069
40 64784 Hannan-Quinn 34.63195
        Excluding the constant, p-value was highest for variable 3 (unem)
\widehat{\text{logcrime}} = -0.709735 - 0.00456 \text{ unem} + 0.000144915 \text{ offi} + 0.8640 \text{ logpop} \\ \xrightarrow{(0.80719)} (0.00903) + (6.1543e-005) (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + (0.0663) + 
                                                                                                                                             (6.1543e-005)
                                    N = 92 \bar{R}^2 = 0.8579 F(3,88) = 184.20 \hat{\sigma} = 0.27968
                                                                              (standard errors in parentheses)
```

First, notice how the coefficients are very different from the one obtain with no logarithm transformation. Here the interpretation is different.  $\beta_2$  is interpreted as the elasticity of crime with respect to pop:

$$\beta_2 = \frac{\Delta \text{crime/crime}}{\Delta \text{pop/pop}}.$$
(5.8)

#### 5.3 Quadratic terms

A one percentage increase in pop will increase crime by 0.864 percent.  $\Delta$  crime/crime is interpreted as a percentage change in crime. For  $\beta_4$  we have:

$$\beta_4 = \frac{\Delta \text{crime/crime}}{\Delta \text{offi}}.$$
(5.9)

Here, a one unit increase in offi is associated with a 0.014% (0.00014  $\times$  100 percent) increase in crime.

## 5.3 Quadratic terms

So far we have bee estimating the marginal effects ( $\beta$ s) that are constant across all possible values of *X*. The simplest way to introduce nonlinearities in the marginal effect is to estimate the model with quadratic terms. For example, let the model be:

$$Y_i = \beta_1 + \beta_2 X_i + \beta_3 X_i^2 + \varepsilon_i. \tag{5.10}$$

In this case the marginal effect of *X* on *Y* is given by:

$$\frac{\Delta Y}{\Delta X} = \beta_2 + 2 \cdot \beta_3 X_i. \tag{5.11}$$

If we want to estimate the marginal effect of experience of wages and in addition we allow for a nonlinear effect we can estimate:

wage<sub>i</sub> = 
$$\beta_1 + \beta_2$$
exper<sub>i</sub> +  $\beta_3$ expersq<sub>i</sub> +  $\varepsilon_i$ , (5.12)

where wage is average hourly earnings, exper is years of experience and expersq is the number of years of experience squared. The Gletl output is the following:

```
      Model 1: OLS, using observations 1-526

      Dependent variable: wage

      coefficient std. error t-ratio p-value

      const 3.72541 0.345939 10.77 1.46e-024 ***

      exper 0.298100 0.0409655 7.277 1.26e-012 ***

      expersq -0.00612989 0.000902517 -6.792 3.02e-011 ***

      Mean dependent var 5.896103 S.D. dependent var 3.693086

      Sum squared resid 6496.147 S.E. of regression 3.524334

      R-squared 0.092769 Adjusted R-squared 0.089300

      F(2, 523)
      26.73982 P-value(F)

      Log-likelihood
      -1407.455 Akaike criterion 2820.910

      Schwarz criterion 2833.706 Hannan-Quinn 2825.920
```

#### 5 Transformations of Variables and Interactions



Fig. 5.1 Predicted values for Equation 5.12

$$\widehat{\text{wage}} = \underbrace{3.72541}_{(0.34594)} + \underbrace{0.298100}_{(0.040966)} \exp \left(-\frac{0.00612989}{(0.00090252)}\right) \exp \left(-\frac{1}{2}\right) \exp \left(-\frac{1}{2$$

Here, the marginal effect of experience on average hourly wage is:

$$\frac{\Delta \text{wage}}{\Delta \text{exper}} = 0.2981 + 2 \cdot (-0.006) \text{exper}$$
$$= 0.2981 - 0.012 \text{exper}.$$

For a person with 2 years of experience, the effect of an additional year of experience on wage is 0.2741 (=0.2981 - 0.012 × 2) and for a person with 15 years of experience, the marginal effect of an additional year of experience is 0.1181 (=0.2981 - $0.012 \times 15$ ). Hence, we can say that for a reasonable range of years of experience, experience has a positive effect on wage. In addition, this effect is smaller as you accumulate more experience.

Figure 5.1 show the fitted regression line along with the 95% confidence interval for the fitted values and the actual data. This figure clearly shows the nonlinear marginal effect and innlustrates how wages increase with experience for about the first 25 years, but then wages decrease later on.

#### 5.4 Interaction terms

## **5.4 Interaction terms**

A second popular approach to allow for the marginal effect to change over different values of *X* is to include interaction terms in the regression equation. For example,

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 (X_2 \times X_3) + \varepsilon_i.$$
(5.13)

In this case the marginal effect of  $X_2$  on Y depends on  $X_3$  is given by:

$$\frac{\Delta Y}{\Delta X_2} = \beta_2 + \beta_3 X_3. \tag{5.14}$$

Consider the next example with the interaction between exper and educ in a wage equation:

wage<sub>i</sub> = 
$$\beta_1 + \beta_2$$
exper<sub>i</sub> +  $\beta_3$ (exper<sub>i</sub> × educ<sub>i</sub>) +  $\varepsilon_i$ , (5.15)

where the marginal effect of experience on wage depends on the level of education:

$$\frac{\Delta \text{wage}}{\Delta \text{exper}} = \beta_2 + \beta_3 \text{educ.}$$
(5.16)

When estimating this equation in Gretl we have to make sure we generate the interaction term first. That is, go to  $Add \rightarrow Define new variable and type:$ 

expereduc = exper\*educ

Then we are ready to estimate the equation via OLS. The regression output is:

Model 1: OLS, using observations 1-526 Dependent variable: wage

coef	ficient	std.	erro	or t-rat	io	p-value	
const 4.8	 8993	0.242	2730	20.15		3.11e-067	***
exper -0.1	88124	0.025	53904	-7.40	9	5.13e-013	***
expereduc 0.0	207731	0.002	21762	25 9.54	5	5.17e-020	***
Mean dependent var	5.8961	03 5	S.D.	dependent	var	3.693086	5
R-squared resid	0.1592	23 <i>1</i>	Adjus	sted R-squ	ared	0.156008	3
F(2, 523)	49.521	75 I	-val	ue(F)		2.01e-20	)
Log-likelihood	-1387.4	49 <i>I</i>	Akaik	e criteri	on	2780.897	7
Schwarz criterion	2793.6	93 H	Hanna	n-Quinn		2785.908	3

```
\widehat{wage} = \underbrace{4.88993}_{(0.24273)} - \underbrace{0.188124}_{(0.025390)} \exp\left(-\frac{10.0207731}{(0.0021762)}\right) \exp\left(-\frac{10.0207731}{(0.0021762)}\right)N = 526 \quad \overline{R}^2 = 0.1560 \quad F(2,523) = 49.522 \quad \widehat{\sigma} = 3.3928
```

(standard errors in parentheses)

Here, the marginal effect of experience on wage is:

$$rac{\Delta ext{wage}}{\Delta ext{exper}} = -0.1881 + 0.0208 ext{educ}$$

For a person with twelve years education (high school), the marginal effect from an additional year of education is 0.0615 (=- $0.1881+0.208\times12$ ). However, with more education the marginal effect is larger. A person with 16 years of education (high school + college) will have a marginal effect of 0.1447 (=- $0.1881+0.208\times16$ ). Notice that for an important range of education the marginal effect is positive, meaning that more experience leads to higher wages. In addition, the effect if larger if you have more education. This means that going to school is not only good because it directly increases your expected wage but also makes additional years of experience more valuable.

# Chapter 6 Analysis with Qualitative Information: Dummy Variables

In previous chapters, the dependent and the independent variables in our regression equations had a *quantitative* meaning. That is, the magnitude of the variable had a useful information, for example, years of education, years of experience, unemployment rate, or wage. In this chapter we will analyze how to introduce *qualitative* information into a regression equation. Example of qualitative information includes marital status, gender, race, industry (manufacturing, retail, etc.) or geographical region (south, north, west, etc.).

# 6.1 Describing qualitative information

Qualitative factors often come in the form of binary information: a person is female of male; a person does or does not own a computer; a person is married or not. In all these cases the relevant information can be captured by a binary variable, also called a dummy variable or zero-one variable. In defining a dummy variable we must decide which event is assigned a value of one and which a value of zero. Table 6.1 shows how two dummy variables (female and married) look in the data set.

person	wage	educ	exper	female	married
1	3.10	11	2	1	0
2	3.24	12	22	1	1
3	3.00	11	2	0	0
4	6.00	8	44	0	1
5	5.30	12	7	0	1
÷	:	÷	÷	÷	:
525	11.56	16	5	0	1
526	3.50	14	5	1	0

Table 6.1 A partial Listing of the Data in Wage.xls



Fig. 6.1 Graph of wage =  $\beta_0 + \delta_0$  female +  $\beta_1$  educ for  $\delta_0 < 0$ .

# 6.2 A single dummy independent variable

The simplest case is when we have a single dummy independent variable. Let's consider the following model:

wage = 
$$\beta_0 + \delta_0$$
 female +  $\beta_1$  educ +  $\varepsilon$  (6.1)

We use the parameter  $\delta_0$  to emphasize the fact that female corresponds to a dummy variable. If the person is a female we have female = 1, and if the person is a male, we have female = 0. The parameter  $\delta_0$  has the following interpretation:  $\delta_0$  is the difference in hourly wage between females and males, given the same amount of education (and the error term  $\varepsilon$ ). Thus, the coefficient  $\delta_0$  determines whether there is discrimination against women: if  $\delta_0 < 0$ , it means that on average, women earn less than men.

The interpretation of  $\delta_0$  (when  $\delta < 0$ ) can be depicted graphically in Figure 6.1 as an *intercept shift* between males an females.

Let's estimate the following more interesting model:

wage = 
$$\beta_0 + \delta_0$$
 female +  $\beta_1$  educ +  $\beta_2$  exper +  $\beta_1$  tenure +  $\varepsilon$  (6.2)

The regression output in Gretl is:

```
Model 2: OLS, using observations 1-526 Dependent variable: wage
```

coeff	icient	std.	erro	r t-ratio	o b	-value	
const -1.56	 794	0.72	4551	-2.164	0.	0309	* *
female -1.81	085	0.26	4825	-6.838	2.	26e-011	***
educ 0.57	1505	0.04	93373	11.58	9.	09e-028	***
exper 0.02	53959	0.01	15694	2.195	Ο.	0286	**
tenure 0.14	1005	0.02	11617	6.663	6.	83e-011	***
Mean dependent var	5.896	103	S.D.	dependent	var	3.69308	36
Sum squared resid	4557.	308	S.E.	of regress	sion	2.95757	72
R-squared	0.363	541	Adju	sted R-squa	ared	0.35865	55
F(4, 521)	74.39	801	P-va	lue(F)		7.30e-5	50
Log-likelihood	-1314.	228	Akai	ke criterio	on	2638.45	55
Schwarz criterion	2659.	782	Hanna	an-Quinn		2646.80	)5

$$\begin{split} \widehat{\text{wage}} &= -1.56794 - 1.81085 \, \text{female} + 0.571505 \, \text{educ} + 0.0253959 \, \text{exper} \\ & (0.72455) & (0.26483) & (0.049337) & (0.011569) \\ & + 0.141005 \, \text{tenure} \\ & (0.021162) & \\ & N = 526 \quad \bar{R}^2 = 0.3587 \quad F(4,521) = 74.398 \quad \hat{\sigma} = 2.9576 \\ & (\text{standard errors in parentheses}) \end{split}$$

Where it is easy to see that  $\delta_0 = -1.81$ . If we want to test the null hypothesis that there is no difference between men and women,  $H_0 : \delta_0 = 0$ . The alternative hypothesis is that there is discrimination against women,  $H_1 : \delta_0 < 0$ . Based on the p-value we reject the null and conclude that there is discrimination, females make two dollars and twenty seven cents less per hour than males. This is after controlling for differences in education, experience and tenure.

It is illustrative to additionally estimate the following equation:

wage = 
$$\beta_0 + \delta_0$$
 female +  $\varepsilon$  (6.3)

where we do not control for education, experience or tenure. The regression output is:

 Model 3: OLS, using observations 1-526

 Dependent variable: wage

 coefficient std. error t-ratio p-value

 const 7.09949 0.210008 33.81 8.97e-134 \*\*\*

 female -2.51183 0.303409 -8.279 1.04e-015 \*\*\*

 Mean dependent var 5.896103 S.D. dependent var 3.693086

 Sum squared resid 6332.194 S.E. of regression 3.476254

 R-squared 0.115667 Adjusted R-squared 0.113979

 F(1, 524)
 68.53668 P-value(F)

 Log-likelihood -1400.732 Akaike criterion 2805.464

 Schwarz criterion 2813.995 Hannan-Quinn 2808.804

6 Analysis with Qualitative Information: Dummy Variables

$$\widehat{\text{wage}} = \frac{7.09949 - 2.51183}_{(0.21001)} \text{ female}$$

$$N = 526 \quad \overline{R}^2 = 0.1140 \quad F(1,524) = 68.537 \quad \widehat{\sigma} = 3.4763$$
(standard errors in parentheses)

The expected (predicted) wage for females is  $\widehat{wage} = 7.099 - 2.5121 = 4.587$ , while the expected wage for males is  $\widehat{wage} = 7.099 - 2.5120 = 7.099$ . This is not controlling for differences in education, experience or tenure. Once we control for those differences, the wage gap between these two groups is smaller and equal to  $\delta_0 = -1.81$ .

What is the interpretation of the coefficient on a dummy variable if the dependent variable is in logs? Here the coefficient has a *percentage* interpretation. Let's say we want to estimate the following equation:

$$\log wage = \beta_0 + \delta_0 \text{female} + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{tenure} + \varepsilon \qquad (6.4)$$

that has the following Gretl estimation output:

$$\begin{split} \log \tilde{wage} &= 0.501348 - 0.301146 \, \text{female} + 0.0874623 \, \text{educ} + 0.00462938 \, \text{exper} \\ & (0.001090) & (0.037246) & (0.0069389) & (0.0016271) \\ & + 0.0173670 \, \text{tenure} \\ & (0.0029762) & \\ & N = 526 \quad \bar{R}^2 = 0.3876 \quad F(4,521) = 84.072 \quad \hat{\sigma} = 0.41596 \\ & (\text{standard errors in parentheses}) \end{split}$$

The coefficient on female,  $\delta_0$ , implies that for the same levels of education, experience, and tenure, women earn approximately 100(0.301) = 30.1% less than men.

# 6.3 Dummy variables for multiple categories

One can use several dummy variables in the same equation. For example, we can add the dummy variable married to Equation 6.3 to obtain:

wage = 
$$\beta_0 + \delta_0$$
 female +  $\delta_1$  married +  $\varepsilon$  (6.5)

In Gretl we have,

$$\widehat{\text{wage}} = \underbrace{\begin{array}{l} 6.18043 - 2.29440 \text{ female} + 1.33948 \text{ married} \\ (0.29634) & (0.30261) & (0.30971) \end{array}}_{N = 526 \quad \overline{R}^2 = 0.1429 \quad F(2,523) = 44.779 \quad \widehat{\sigma} = 3.4190 \\ & (\text{standard errors in parentheses}) \end{array}}$$

The coefficient on married gives the (approximate) difference in wages between married and non married individuals. Based on these results, married individuals

6.4 Incorporating ordinal information

have higher hourly wages. On important restriction in Equation 6.5 is that it restricts the effect of marital status on wages is the same whether you are male of female. If we are interested in this difference we can estimate an alternative model with additional categories. In particular we need four categories: (1) married men, (2) married women (3) single men, and (4) single woman. We must select a base group (for example, single men) and create the dummy variables for the other three groups.

The equation we want to estimate is:

$$\log wage = \beta_0 + \delta_0 marrmale + \delta_1 marrfem + \delta_2 singfem + \varepsilon$$
 (6.6)

and the estimation output is:

$$\widehat{\log wage} = \underbrace{1.5201}_{(0.050987)} + \underbrace{0.4267}_{(0.061554)} \text{ marrmale} - \underbrace{0.0797}_{(0.065524)} \text{ marrfem} - \underbrace{0.1316}_{(0.066804)} \text{ singfem}$$

$$N = 526 \quad \overline{R}^2 = 0.2087 \quad F(3, 522) = 47.149 \quad \widehat{\sigma} = 0.47284$$
(standard errors in parentheses)

The interpretation of each of the  $\delta$  coefficients is with respect to the base group. For example  $\delta_2 = 0.1316$  means that single females earn approximately 13.16% lower hourly wages than single men (the base group).

# 6.4 Incorporating ordinal information

Suppose we want to estimate the effect of city credit ratings on the municipal bond interest rate (MBR). The credit rating (CR) is an ordinal variable and suppose it goes from zero (worst credit) to four (best credit). Under these consideration, a potential candidate for our model is:

$$MBR = \beta_0 + \beta_1 CR + other factors + \varepsilon$$
(6.7)

where *other factors* are just other variables in the model. On concern with this specification is that it is hard to interpret one unit increase in CR. It is easy to talk about an additional year of education or an additional year of income, but credit ratings usually have only an ordinal meaning. Moreover, it is restrictive to assume that each additional unit increase in CR has the same effect on MBR. An alternative approach is to create separate dummy variables for each of the values of CR, that is,

6 Analysis with Qualitative Information: Dummy Variables

$$CR_1 = 1 \text{ if } CR = 1$$
  
= 0 otherwise.  
$$CR_2 = 1 \text{ if } CR = 2$$
  
= 0 otherwise.  
$$CR_3 = 1 \text{ if } CR = 3$$
  
= 0 otherwise.  
$$CR_4 = 1 \text{ if } CR = 4$$
  
= 0 otherwise.

Then we can focus on estimating the following model:

$$MBR = \beta_0 + \delta_1 CR_1 + \delta_2 CR_2 + \delta_3 CR_3 + other factors + \varepsilon$$
(6.8)

Again, we omit one category (CR<sub>4</sub>) and the interpretation of the dummy coefficients is relative to the omitted category. For example,  $\delta_2$  represents the difference in municipal bond interest rate between ratings CR<sub>2</sub> and CR<sub>4</sub>.

# 6.5 Interactions involving dummy variables

Just as quantitative variables can have interactions, so can dummy variables. Actually, we revisit the estimation of Equation 6.6 to see that the same results can be achieved by including the interaction term between female and married. The model we want to estimate is:

 $\log wage = \beta_0 + \delta_0 \text{female} + \delta_1 \text{married} + \delta_2 (\text{female} \times \text{married}) + \varepsilon$  (6.9)

Estimating in Gretl we have:

$$\log wage = \frac{1.5201}{(0.050987)} - \frac{0.1316}{(0.066804)} \text{ female} + \frac{0.4267}{(0.061554)} \text{ married}$$

$$- \frac{0.3748}{(0.085708)} \text{ female} \times \text{married}$$

$$N = 526 \quad \bar{R}^2 = 0.2087 \quad F(3,522) = 47.149 \quad \hat{\sigma} = 0.47284$$
(standard errors in parentheses)

Notice that this regression output is equivalent as the one obtained from Equation 6.6.



 $\label{eq:Fig.6.2} \textbf{Fig. 6.2 Graph of wage} = \beta_0 + \delta_0 \texttt{female} + \beta_1 \texttt{educ} + \delta_1 \texttt{educ} \times \texttt{female}.$ 

# 6.5.1 Allowing for different slopes

Consider the case where we want to estimate the effect of education on hourly wage and in addition, we want for the marginal effect to change based on your gender. This can be done by interacting the educ with female and estimating the following model:

wage = 
$$\beta_0 + \delta_0$$
 female +  $\beta_1$  educ +  $\delta_1$  (female × educ) +  $\varepsilon$  (6.10)

A graphical approach to this problem in presented in Figure 6.2. The output in Gretl is

$$\widehat{\text{wage}} = \underbrace{0.200496}_{(0.84356)} - \underbrace{1.19852}_{(1.3250)} \underbrace{\text{female}}_{(0.064223)} + \underbrace{0.064223}_{(0.064223)} - \underbrace{0.0859990}_{(0.10364)} \underbrace{\text{female}}_{R} \times educ$$

$$T = 526 \quad \overline{R}^2 = 0.2555 \quad F(3,522) = 61.070 \quad \widehat{\sigma} = 3.1865$$
(standard errors in parentheses)

### 6.5.2 Testing for differences in regression functions across groups

So far we saw that interacting a dummy variable with other independent variables is a powerful tool. Now, we can use this tool to test the null hypothesis that two groups follow the same regression function, against the alternative that one or more of the slopes differs across the two groups. Suppose we want to test whether the same regression model describe college GPA for males and for females. The model is

$$cumgpa = \beta_0 + \beta_1 sat + \beta_2 hsperc + \beta_3 tothrs + \varepsilon, \qquad (6.11)$$

where cumgpa is cumulative college GPA, sat is the SAT score, hsperc is the high school rank percentile, and tothrs is the total hours of college courses. The regression results in Gretl are

$$\begin{split} & \text{cumgpa} = \underbrace{0.929111}_{(0.22855)} + \underbrace{0.0009028}_{(0.000208)} \text{ sat} - \underbrace{0.006379}_{(0.00157)} \text{ hsperc} + \underbrace{0.01198}_{(0.000931)} \text{ tothrs} \\ & N = 732 \quad \bar{R}^2 = 0.2323 \quad F(3,728) = 74.717 \quad \hat{\sigma} = 0.86711 \\ & \text{(standard errors in parentheses)} \end{split}$$

To allow for a difference in the intercept we just need to include the dummy variable female. Then, to allow differences in the slope parameters we need to include interaction terms for each of the variables and female. That is

$$\begin{array}{ll} \text{cumgpa} = \beta_0 + \delta_0 \text{female} + \beta_1 \text{sat} + \delta_1 \text{sat} \cdot \text{female} & (6.12) \\ & + \beta_2 \text{hsperc} + \delta_2 \text{hsperc} \cdot \text{female} \\ & + \beta_3 \text{tothrs} + \delta_3 \text{tothrs} \cdot \text{female} + \varepsilon \end{array}$$

The parameter  $\delta_0$  is the difference in the intercepts between females and males,  $\delta_1$  is the slope difference with respect to sat between females and males, and so on. The null hypothesis that cumppa follows the same model for females and males is

$$H_0: \delta_0 = \delta_1 = \delta_2 = \delta_3 = 0 \tag{6.13}$$

If at least one of the  $\delta_j$  is different from zero, then the model is different for men and women. After creating the interaction terms, the estimated model in Gretl is

```
Model 2: OLS, using observations 1-732 Dependent variable: cumgpa
```

	coefficient	std. error	t-ratio	p-value	
const	1.21398	0.264828	4.584	5.37e-06	***
sat	0.000611312	0.000235026	2.601	0.0095	***
hsperc	-0.00596745	0.00177646	-3.359	0.0008	***
tothrs	0.0103004	0.00109284	9.425	5.65e-020	***
female	-1.11364	0.528539	-2.107	0.0355	**
satfemale	0.00111674	0.000500034	2.233	0.0258	**
hspercfemale	5.07597e-05	0.00410253	0.01237	0.9901	
tothrsfemale	0.00555989	0.00206958	2.686	0.0074	***

2.080861	S.D. dependent var	0.989617
534.3092	S.E. of regression	0.859067
0.253652	Adjusted R-squared	0.246436
35.15106	P-value(F)	2.54e-42
-923.4440	Akaike criterion	1862.888
1899.654	Hannan-Quinn	1877.071
	2.080861 534.3092 0.253652 35.15106 -923.4440 1899.654	2.080861 S.D. dependent var 534.3092 S.E. of regression 0.253652 Adjusted R-squared 35.15106 P-value(F) -923.4440 Akaike criterion 1899.654 Hannan-Quinn

```
\begin{split} \widehat{\text{cumgpa}} &= \underbrace{1.21398}_{(0.26483)} + \underbrace{0.000611312}_{(0.00023503)} \text{ sat} - \underbrace{0.00596745}_{(0.0017765)} \text{ sperc} + \underbrace{0.0103004}_{(0.0010928)} \text{ tothrs} \\ &- \underbrace{1.11364}_{(0.52854)} \text{ female} + \underbrace{0.00111674}_{(0.00050003)} \text{ satfemale} + \underbrace{5.07597e}_{(0.0041025)} \text{ spercfemale} \\ &- \underbrace{0.00555989}_{(0.0020696)} \text{ tothrsfemale} \\ &- \underbrace{0.0020696}_{(0.0020696)} \text{ satfemale} + \underbrace{0.00111674}_{(0.0020696)} \text{ satfemale} + \underbrace{0.0011025}_{(0.0011025)} \text{ satfemale} \\ &- \underbrace{0.0020696}_{(0.0020696)} \text{ satfemale} + \underbrace{0.00111674}_{(0.0020696)} \text{ satfemale} + \underbrace{0.0011025}_{(0.0011025)} \text{ satfemale} \\ &- \underbrace{0.0020696}_{(0.0020696)} \text{ satfemale} + \underbrace{0.00111674}_{(0.0020696)} \text{ satfemale} + \underbrace{0.0011025}_{(0.0020696)} \text{ satfemale} \\ &- \underbrace{0.0020696}_{(0.0020696)} \text{ satfemale} + \underbrace{0.00111674}_{(0.0020696)} \text{ satfemale} + \underbrace{0.0011025}_{(0.0020696)} \text{ satfemale} \\ &- \underbrace{0.0020696}_{(0.0020696)} \text{ satfemale} + \underbrace{0.00111674}_{(0.0020696)} \text{ satfemale} + \underbrace{0.0011025}_{(0.0020696)} \text{ satfemale} \\ &- \underbrace{0.0020696}_{(0.0020696)} \text{ satfemale} + \underbrace{0.00111674}_{(0.0020696)} \text{ satfemale} + \underbrace{0.0011025}_{(0.0020696)} \text{ satfemale} \\ &- \underbrace{0.0020696}_{(0.0020696)} \text{ satfemale} + \underbrace{0.00111674}_{(0.0020696)} \text{ satfemale} + \underbrace{0.0011025}_{(0.0020696)} \text{ satfemale} + \underbrace{0.0020696}_{(0.0020696)} \text{ satfemale} + \underbrace{0.0020696}_{(0.0020696)} \text{ satter} + \underbrace{0.0020696}_{(0.0020696}_{(0.0020696)} \text{ satter} + \underbrace{0.0020696}_{(0.0020696}_{(0.0020696)}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(0.0020696}_{(
```

N = 732  $\bar{R}^2 = 0.2464$  F(7,724) = 35.151  $\hat{\sigma} = 0.85907$ (standard errors in parentheses)

Now, to test the null hypothesis presented in Equation 6.13 from the window that shows the regression output, we need to go to Tests  $\rightarrow$  Omit variables and a new window will open. We then have to select the variables to omit. There are female, satfemale, hspercfemale, and tothrsfemale. This will estimate the restricted model and the comparison between the restricted model (Equation 6.11) and the full model (Equation 6.12),

Model 3: OLS, using observations 1-732 Dependent variable: cumgpa

	coeffi	cient	std	. error	t-ratio	p-value	
const	0.929	 111	0.2	 28552	4.065	5.32e-05	* * *
sat	0.000	902834	0.0	0207870	4.343	1.60e-05	* * *
hsperc	-0.006	37913	0.0	0156785	-4.069	5.24e-05	* * *
tothrs	0.011	9779	0.0	00931383	12.86	2.96e-034	* * *
Mean depend	ent var	2.0808	861	S.D. dep	endent var	0.989617	
R-squared	resta	0.2354	16	Adjusted	R-squared	0.232265	
F(3, 728)		/4./1/	07	P-value(	F')	3.87e-42	
Log-likelih	.ood	-932.27	97	Akaike c	riterion	1872.559	
Schwarz cri	terion	1890.9	942	Hannan-Q	uinn	1879.651	
Comparison	of Model	2 and M	Iodel	3:			
Null hypo female,	thesis:	the regr le, hspe	ressio	on parame male, tot	ters are ze hrsfemale	ro for the	variable

Test statistic: F(4, 724) = 4.4227, with p-value = 0.00154347 Of the 3 model selection statistics, 1 has improved.

The *F* statistics that Gretl is reporting comes from

6 Analysis with Qualitative Information: Dummy Variables

$$F = \frac{RSS - RSS_{\rm UR}}{RSS_{\rm UR}} \cdot \frac{n - 2k}{q},\tag{6.14}$$

where RSS is the residual sum of squares of the model estimates in Equation 6.11 and  $RSS_{UR}$  is the unrestricted model in Equation 6.12. *n* is the sample size, *k* is the number of parameters we are estimating, and *q* is the number of restrictions when comparing the model in Equation 6.11 and in Equation 6.12. Substituting the values we obtain,

$$F = \frac{547.3649 - 534.3092}{534.3092} \cdot \frac{732 - 2 \cdot 4}{4} = 0.024434 \cdot 181 = 4.4227, \tag{6.15}$$

An alternative way to calculate this F statistic is to follow the formula,

$$F = \frac{RSS - (RSS_1 + RSS_2)}{RSS_1 + RSS_2} \cdot \frac{n - 2k}{k},$$
(6.16)

where RSS is the residual sum of squares of the model estimates in Equation 6.11.  $RSS_1$  and  $RSS_2$  are the residual sum of squares of the model estimated in Equation 6.11 using only the females in the sample ( $RSS_1$ ) and using only the males in the sample ( $RSS_2$ ). As before, *n* is the sample size and *k* is the number of parameters we are estimating. The estimation of Equation 6.11 with just females is:

```
Model 5: OLS, using observations 1-180 Dependent variable: cumgpa
```

	coeffi	cient	std.	error	t-ratio	p-value	
const	0.100	 346	0.48	 1095	0.2086	0.8350	
sat	0.001	72805	0.00	046421	6 3.723	0.0003	***
hsperc	-0.005	91669	0.00	388949	-1.521	0.1300	
tothrs	0.015	3603	0.00	184854	8.580	4.82e-015	***
Mean depende	nt var	2.268	611	S.D. (	dependent var	1.126549	)
Sum squared	resid	143.6	897	S.E.	of regression	0.903559	)
R-squared		0.367	483	Adjust	ted R-squared	0.356702	2
F(3, 176)		34.08	447	P-val:	ue(F)	2.03e-17	7
Log-likeliho	od	-235.1	319	Akaik	e criterion	478.2638	3
Schwarz crit	erion	491.0	356	Hannaı	n-Ouinn	483.4422	2

#### and with just males is:

Model 6: OLS, using observations 1-552 Dependent variable: cumgpa

	coefficient	std. error	t-ratio	p-value	
const	1.21398	0.260270	4.664	3.90e-06	***
sat	0.000611312	0.000230981	2.647	0.0084	***
hsperc	-0.00596745	0.00174588	-3.418	0.0007	***
tothrs	0.0103004	0.00107403	9.590	3.06e-020	***
Mean depend	ent var 2.019	638 S.D. dep	endent var	0.933655	

#### 6.6 The dummy variable trap

Sum squared resid	390.6194	S.E. of regression	0.844280
R-squared	0.186740	Adjusted R-squared	0.182288
F(3, 548)	41.94377	P-value(F)	2.06e-24
Log-likelihood	-687.8093	Akaike criterion	1383.619
Schwarz criterion	1400.873	Hannan-Quinn	1390.360

Using the formula in Equation 6.17,

$$F = \frac{547.3649 - (143.6897 + 390.6194)}{143.6897 + 390.6194} \cdot \frac{732 - 2 \cdot 4}{4} = 0.024434 \cdot 181 = 4.4227,$$
(6.17)

which is the same result as in Equation 6.15. This version of the F test is know also as the *Cho* test. A large F statistic is evidence against the null hypothesis. In our example the F statistic of 4.4227 has an associated p-value of 0.0015, below the usual 0.05 (or 5%). Hence, we reject the null hypothesis that there is no difference between the equation for females and the equation for males. This means that there is difference and we are better off estimating Equation 6.12 instead of Equation 6.11.

The key to estimate Equation 6.11 with just the female portion of the data change the sample. To do this go to Sample  $\rightarrow$  Restrict, based on criterion..., then after a new window shows up, select the "use dummy variable" and then female. Once the sample is restricted, just estimate the model using Ordinary Least Squares again.

## 6.6 The dummy variable trap

The dummy variable trap occurs when there is an exact linear relationship among the variables in the regression model. That is the reason why we do not include female and male in the same regression equation because female + male = 1. The same occurs when we have more than one category and we should always omit one of the categories (base group). Than is why singmen does not appear in Equation 6.6 (marrmale + singmale + marrfem + singfem = 1).

# Chapter 7 Specification of Regression Variables

So far we assumed we know what are the variables that needed to be in our regression model. However, what happens if we include in the regression model a variable that should not be there? What if we leave out a variable that should be included? Can we a proxy for a variable that we do not observe? These are the main question this chapter will address.

# 7.1 Model specification

What happens in practice is that it is difficult to be sure about the correct specification of the regression model. While theory may help, it usually depends on simplifying assumptions that may not necessarily hold. The properties of the regression estimates depend crucially on the validity of the specification of the model. The following is a quick summary of the consequences of misspecifying the regression model:

- 1. If you leave out a variable that should be included. The regression estimates are potentially biased. The standard errors of the coefficients and the corresponding t and F tests are in general invalid.
- 2. If you include a variable that should not be in the model. The coefficients will not be biased, but they are potentially inefficient.

## 7.2 Omitting a variable

## 7.2.1 The bias problem

Suppose that the true regression model that we should be estimated is given by

7 Specification of Regression Variables

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + u. \tag{7.1}$$

However, we do not have the variable  $X_3$  or maybe we have it but we do not include it in the model. Hence, we estimate the following model

$$Y = \beta_1 + \beta_2 X_2 + u. \tag{7.2}$$

Then the predicted or fitted values are

$$\hat{Y} = b_1 + b_2 X_2 \tag{7.3}$$

Recall from previous chapters that the formula to estimate  $b_2$  is given by

$$b_2 = \frac{\sum_{i=1}^n (X_{2i} - \bar{X}_2)(Y_i - \bar{Y})}{\sum_{i=1}^n (X_{2i} - \bar{X}_2)^2}$$
(7.4)

We say that  $b_2$  is unbiased if its expected value is equal to the true population parameter  $\beta_2$ . If we plug Equation 7.1 into Equation 7.4 and take expectations we obtain

$$E[b_2] = E\left[\frac{\sum_{i=1}^n (X_{2i} - \bar{X}_2)(Y_i - \bar{Y})}{\sum_{i=1}^n (X_{2i} - \bar{X}_2)^2}\right]$$
  
=  $\beta_2 + \beta_3 \frac{\sum_{i=1}^n (X_{2i} - \bar{X}_2)(X_{3i} - \bar{X}_3)}{\sum_{i=1}^n (X_{2i} - \bar{X}_2)^2}.$  (7.5)

For  $b_2$  to be unbiased we need that the second term on the right-hand side be equal to zero. This term is known as the *omitted variable bias* and it will be zero if  $\beta_3 = 0$  or if  $\sum_{i=1}^{n} (X_{2i} - \bar{X}_2)(X_{3i} - \bar{X}_3) / \sum_{i=1}^{n} (X_{2i} - \bar{X}_2)^2$  is equal to zero. Then the conditions for  $b_2$  to be unbiased in the estimation of Equation 7.2 are:

- 1. That  $X_3$  does not affect *Y*. That is,  $\beta_3 = 0$ .
- 2. That  $X_2$  and  $X_3$  are linearly uncorrelated. That is, the slope coefficient when we regress  $X_3$  on  $X_2$  os zero,  $\sum_{i=1}^{n} (X_{2i} \bar{X}_2)(X_{3i} \bar{X}_3)]/[\sum_{i=1}^{n} (X_{2i} \bar{X}_2)^2 = 0.$

# 7.2.2 Invalid statistical tests

When a variable is omitted from the model, the standard errors of the coefficients and the texts statistics are generally invalid. This means that the t and F tests cannot be used.

#### 7.2.3 *Example*

Consider the case where the true model to explain wages is given by

#### 7.2 Omitting a variable

wage = 
$$\beta_1 + \beta_2$$
educ +  $\beta_3$ ability +  $u$ . (7.6)

That is, your wage is determined by your number of years of formal education (educ) and your ability. The problem in this equation is that actually it is very difficult to measure ability. Hence, we decide omit it and estimate the following model

wage = 
$$\beta_1 + \beta_2$$
educ +  $u$ . (7.7)

What is the problem with the estimate of  $\beta_2$  if we use Equation 7.7? Is is biased! To get an idea of the size of the bias we will proxy ability with another variable, IQ. Equation 7.5 becomes

$$E[b_2] = \beta_2 + \beta_3 \frac{\sum_{i=1}^{n} (\text{educ}_i - \overline{\text{educ}})(IQ_i - \overline{IQ})}{\sum_{i=1}^{n} (\text{educ}_i - \overline{\text{educ}})^2}.$$
 (7.8)

Notice that we can actually analyze if the bias is positive or negative based on the signs of the second part on the right-hand size. It seems that  $\beta_3$  should be positive because higher ability (or IQ) should be correlated positively with wages. Moreover, the part that multiplies  $\beta_3$  should also be positive because education and ability (or IQ) seem to be positively correlated. Hence, the whole second part on the right-hand side is positive, implying that  $\beta_2$  is biased upwards. This means that on average we will be getting a larger coefficient (by estimating Equation 7.7) than the true coefficient (if we were estimating the true Equation 7.6).

Let's look at this empirically by estimating Equations 7.6 and 7.7 with real data (where we use IQ in place of ability):

```
Model 1: OLS, using observations 1-935 Dependent variable: wage
```

const       146.952       77.7150       1.891       0.0589       *         educ       60.2143       5.69498       10.57       9.35e-025       *         Mean dependent var       957.9455       S.D. dependent var       404.3608         Sum squared resid       1.36e+08       S.E. of regression       382.3203         R-squared       0.107000       Adjusted R-squared       0.106043         F(1, 933)       111.7929       P-value(F)       9.35e-25         Log-likelihood       -6885       458       Akaike criterion       13774		coeffic	ient s	td. erro	r t-ratio	o p-val	ue
Mean dependent var         957.9455         S.D. dependent var         404.3608           Sum squared resid         1.36e+08         S.E. of regression         382.3203           R-squared         0.107000         Adjusted R-squared         0.106043           F(1, 933)         111.7929         P-value(F)         9.35e-25           Log-likelihood         -6885.458         Akaike criterion         13774.92	const educ	146.952 60.21	 2 43	 77.7150 5.69498	1.891 10.57	0.0589 9.35e-	* 025 ***
Schwarz criterion 13784.60 Hannan-Quinn 13778.61	Mean depende Sum squared R-squared F(1, 933) Log-likeliho Schwarz crit	nt var resid od - erion	957.945 1.36e+0 0.10700 111.792 -6885.45 13784.6	5 S.D. 8 S.E. 0 Adju 9 P-va 8 Akai 0 Hann	dependent of regres: sted R-squa lue(F) ke criteric an-Ouinn	var 404 sion 382 ared 0.1 9.3 on 137 137	.3608 .3203 06043 5e-25 74.92 78.61

$$\widehat{wage} = \frac{146.952 + 60.2143}{(77.715)} \underbrace{+ 60.2143}_{(5.6950)} \underbrace{+ 60.2143$$

N = 935  $\bar{R}^2 = 0.1060$  F(1,933) = 111.79  $\hat{\sigma} = 382.32$ 

(standard errors in parentheses)

Model 2: OLS, using observations 1-935 Dependent variable: wage

coefficient std. error t-ratio p-value

7 Specification of Regression Variables

 $N = 935 \quad \bar{R}^2 = 0.1320 \quad F(2,932) = 72.015 \quad \hat{\sigma} = 376.73$ (standard errors in parentheses)

The empirical results are consistent with our theoretical analysis. The estimate of  $\beta_2$  in Equation 7.7 is too large (upward biased). The bias can be obtained separately by estimating a regression of IQ on educ and then plugging the results in Equation 7.8

Model 3: OLS, using observations 1-935

Dependent variable: IQ coefficient std. error t-ratio p-value const 53.6872 2.62293 20.47 3.36e-077 \*\*\* educ 3.53383 0.192210 18.39 1.16e-064 \*\*\* Mean dependent var 101.2824 S.D. dependent var 15.05264 Sum squared resid 155346.5 S.E. of regression 12.90357 R-squared 0.265943 Adjusted R-squared 0.265157 F(1, 933) 338.0192 P-value(F) 1.16e-64 Log-likelihood -3716.973 Akaike criterion 7437.946 Schwarz criterion 7447.627 Hannan-Quinn 7441.637

> $\widehat{IQ} = 53.6872 + 3.53383 \text{ educ}$ (2.6229) (0.19221)

$$N = 935$$
  $\bar{R}^2 = 0.2652$   $F(1,933) = 338.02$   $\hat{\sigma} = 12.904$ 

(standard errors in parentheses)

Replacing the valued in Equation 7.8

$$E[b_2] = \beta_2 + \beta_3 \frac{\sum_{i=1}^{n} (\text{educ}_i - \text{educ})(IQ_i - \overline{IQ})}{\sum_{i=1}^{n} (\text{educ}_i - \overline{\text{educ}})^2}.$$

$$= \beta_2 + 5.13796 \times 3.53383$$

$$= \beta_2 + 18.15667$$
(7.9)

That is exactly the difference between the coefficients in Equations 7.6 and 7.7, 60.2143 - 42.0576 = 18.15667.

#### 7.3 Including a variable that should not be included

Suppose that the true population model is given by

$$Y = \beta_1 + \beta_2 X_2 + u. \tag{7.10}$$

However, for some season you include  $X_3$  and end up estimating the following model

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + u. \tag{7.11}$$

In a regression model like Equation 7.11 with two variables ( $X_2$  and  $X_3$ ) the OLS estimator for  $b_2$  is given by

$$b_{2} = \frac{\sum_{i=1}^{n} (X_{2i} - \bar{X}_{2})(Y_{i} - \bar{Y}) \sum_{i=1}^{n} (X_{3i} - \bar{X}_{3})^{2}}{\sum_{i=1}^{n} (X_{2i} - \bar{X}_{2})^{2} \sum_{i=1}^{n} (X_{3i} - \bar{X}_{3})^{2} - \left(\sum_{i=1}^{n} (X_{2i} - \bar{X}_{2})(X_{3i} - \bar{X}_{3})\right)^{2}} - \frac{\sum_{i=1}^{n} (X_{3i} - \bar{X}_{3})(Y_{i} - \bar{Y}) \sum_{i=1}^{n} (X_{2i} - \bar{X}_{2})(X_{3i} - \bar{X}_{3})}{\sum_{i=1}^{n} (X_{2i} - \bar{X}_{2})^{2} \sum_{i=1}^{n} (X_{3i} - \bar{X}_{3})^{2} - \left(\sum_{i=1}^{n} (X_{2i} - \bar{X}_{2})(X_{3i} - \bar{X}_{3})\right)^{2}}$$
(7.12)

Which is certainly different from the OLS estimator for  $b_2$  in Equation 7.10,

$$b_2 = \frac{\sum_{i=1}^n (X_{2i} - \bar{X}_2)(Y_i - \bar{Y})}{\sum_{i=1}^n (X_{2i} - \bar{X}_2)^2}$$
(7.13)

Interestingly,  $b_2$  in both Equations (7.12 and 7.13) is unbiased,  $E(b_2) = \beta_2$ . Hence, estimating the effect of  $X_2$  on Y will yield unbiased estimates even if we include irrelevant variables. Then, what is the problem? Including irrelevant variables will inflate the standard errors of the coefficients. This means that the estimate  $b_2$  from Equation 7.11 will be inefficient. The implied population variance of  $b_2$  in Equation 7.11 is

$$\sigma_{b_2}^2 = \frac{\sigma_u^2}{\sum_{i=1}^n (X_{2i} - \bar{X}_2)^2} \cdot \frac{1}{(1 - r_{X_2 X_3}^2)}$$
(7.14)

where  $r_{X_2X_3}^2$  is the correlation coefficient between  $X_2$  and  $X_3$ , while the population variance of  $b_2$  in Equation 7.10 is

$$\sigma_{b_2}^2 = \frac{\sigma_u^2}{\sum_{i=1}^n (X_{2i} - \bar{X}_2)^2}.$$
(7.15)

Notice that because  $0 \le r_{X_2X_3}^2 \le 1$ , the population variance in Equation 7.15 is larger than the implied population variance in Equation 7.14. Actually, they will be equal if  $r_{X_2X_3}^2 = 0$ , that is, if  $X_2$  and  $X_3$  are linearly uncorrelated. Moreover, when linearly un-

correlated  $\sum_{i=1}^{n} (X_{2i} - \bar{X}_2)(X_{3i} - \bar{X}_3) = 0$ , then Equation 7.12 reduces to 7.13, meaning that including  $X_3$  in the equation will not affect the estimation of  $\beta_2$ . While the population variances are the same, the estimated (sample) variances will still differ due to a reduction in the degrees of freedom.

# 7.3.1 Example

Consider the following model where we want to see how age affects the likelihood of being married. Are older people more likely to be married? Well, let's estimate the exact response of married to age,<sup>1</sup>

married = 
$$\beta_1 + \beta_2$$
age +  $\epsilon$  (7.16)

The estimation results from Gretl are

```
Model 1: OLS, using observations 1-935 Dependent variable: married
```

cient	std.	error	t-ratio	p-value	
935	0.107	7608	5.027	5.98e-07	***
6442	0.003	323870	3.287	0.0011	***
0.8930	)48	S.D. dep	pendent var	0.3092	217
88.282	274	S.E. of	regressior	n 0.3076	508
0.0114	145	Adjusted	d R-squared	d 0.0103	385
10.801	60	P-value	(F)	0.0010	)52
-223.40	)66	Akaike	criterion	450.81	L33
460.49	944	Hannan-9	Quinn	454.50	)47
	cient 935 6442 0.8930 88.282 0.0114 10.801 -223.40 460.49	cient std. 935 0.107 6442 0.003 0.893048 88.28274 0.011445 10.80160 -223.4066 460.4944	cient std. error 935 0.107608 6442 0.00323870 0.893048 S.D. dep 88.28274 S.E. of 0.011445 Adjuster 10.80160 P-value -223.4066 Akaike of 460.4944 Hannan-	cient std. error t-ratio 935 0.107608 5.027 6442 0.00323870 3.287 0.893048 S.D. dependent van 88.28274 S.E. of regression 0.011445 Adjusted R-squared 10.80160 P-value(F) -223.4066 Akaike criterion 460.4944 Hannan-Quinn	cient std. error t-ratio p-value 935 0.107608 5.027 5.98e-07 6442 0.00323870 3.287 0.0011 0.893048 S.D. dependent var 0.3092 88.28274 S.E. of regression 0.3076 0.011445 Adjusted R-squared 0.0103 10.80160 P-value(F) 0.0010 -223.4066 Akaike criterion 450.83 460.4944 Hannan-Quinn 454.50

$$\widehat{R^2} = 0.540935 + 0.0106442 \text{ age} \\ (0.10761) + (0.0032387) \\ N = 935 \quad \overline{R}^2 = 0.0104 \quad F(1,933) = 10.802 \quad \widehat{\sigma} = 0.30761 \\ (\text{standard errors in parentheses})$$

If the average age in the sample is 33 years of age, the predicted value for married is 89.2 (married =  $0.5409 + 0.0106 \times 33$ ). This means that if you are 33 years old, the probability that you are married is 89.2%. In addition, every year you get older, the probability that you are married increases by 0.011 or about 1.%. For some reason you think that only fools get married and then you decide to wrongly estimate the model

married = 
$$\beta_1 + \beta_2$$
age +  $\beta_3$ IQ +  $\varepsilon$  (7.17)

<sup>&</sup>lt;sup>1</sup> Because married is actually a dummy variable this is a linear probability model, a type of model that we will see in detail in Chapter 9.

#### 7.4 Testing a linear restriction

where the variable IQ is  $X_3$  in Equation 7.11 and should not be in the model. The estimation results from Gretl are

```
Model 2: OLS, using observations 1-935
Dependent variable: married
                 coefficient std. error t-ratio p-value
   _____
                 0.563197 0.129804 4.339 1.59e-05 ***
0.0106007 0.00324337 3.268 0.0011 ***
   const
  age
  IQ
                -0.000205573 0.000669635 -0.3070 0.7589
Mean dependent var 0.893048 S.D. dependent var 0.309217

        Sum squared resid
        88.27381
        S.E. of regression
        0.307757

        R-squared
        0.011545
        Adjusted R-squared
        0.009424

        F(2, 932)
        5.442677
        P-value(F)
        0.004467

                                                                        0.004467
452.7187
458.2559

      F(2, 932)
      5.442677
      P-value(F)

      Log-likelihood
      -223.3594
      Akaike criterion

Schwarz criterion 467.2404 Hannan-Quinn
         married = 0.563197 + 0.0106007 age - 0.000205573 IQ
                                                         (0.00066963)
                          (0.12980)
                                    (0.0032434)
         N = 935 \bar{R}^2 = 0.0094 F(2,932) = 5.4427 \hat{\sigma} = 0.30776
                         (standard errors in parentheses)
```

Not surprisingly, the effect of IQ on married is not significant. This means that fools are not more likely to be married. However, the results do not necessarily support the conjecture that higher IQ is associated with married people either. Nevertheless, including IQ does not seems to help in the estimation of  $\beta_2$ . As we have seen theoretically, the estimate of the second equation is less efficient as can be appreciated from its larger standard error (0.003243 > 0.003239).

# 7.4 Testing a linear restriction

Testing linear restriction on the regression coefficients is sometimes very useful. Consider the following regression model,

$$\log wage = \beta_1 + \beta_2 exper + \beta_3 educ + \varepsilon$$
(7.18)

The regression output in Gretl is

Model 1: OLS, using observations 1-935 Dependent variable: logwage

	coefficient	std. error	t-ratio	p-value	
const	5.50271	0.112037	49.12	8.13e-261	***
educ	0.0777820	0.00657687	11.83	3.62e-030	***
exper	0.0197768	0.00330251	5.988	3.02e-09	* * *

#### 7 Specification of Regression Variables

Mean dependent var	6.779004	S.D. dependent var	0.421144
Sum squared resid	143.9786	S.E. of regression	0.393044
R-squared	0.130859	Adjusted R-squared	0.128994
F(2, 932)	70.16174	P-value(F)	4.13e-29
Log-likelihood	-452.0704	Akaike criterion	910.1407
Schwarz criterion	924.6624	Hannan-Quinn	915.6779

logwage = 
$$5.50271 + 0.0777820$$
 educ +  $0.0197768$  exper  
(0.11204) + (0.0065769) + (0.0033025) =  $N = 935$   $\bar{R}^2 = 0.1290$   $F(2,932) = 70.162$   $\hat{\sigma} = 0.39304$  (standard errors in parentheses)

Let's say that we want to text whether the effect of a year on education on wages is the same as the effect of a year of experience of wages. That is, we want to text the following null hypothesis,

$$H_0: \beta_2 = \beta_3 \tag{7.19}$$

While it may be tempting to just look and compare the regression estimates  $b_2$  and  $b_3$ , this approach is not correct. Remember that  $b_2$  and  $b_3$  are just estimates and are not the unknown  $\beta_2$  and  $\beta_3$ . The statistically correct approach is to run an auxiliary restricted regression where we force  $b_2 = b_3$ . Then, we have to compare if the regression fit with the *restricted* coefficients is significantly lower that the regression fit with the *unrestricted* (original) regression. To do this we calculate the residual sum of squares from the restricted model ( $RSS_U$ ) and calculate the following F statistic:

$$F_{r,n-k} = \frac{(RSS_R - RSS_U)/r}{RSS_U/(n-k)}$$
(7.20)

where *F* is distributed with *r* and n - k degrees of freedom. The number of restrictions *r* is equal to one in our example.

This is done automatically in Gretl. After you estimate the unrestricted regression model, in the regression output window you have to go to Tests  $\rightarrow$  Linear restrictions and a new window will open. In the new window you have to type the command b[educ] - b[exper] = 0 to obtain

const	6.24122	0.0877816	71.10	0.0000	***				
educ	0.0214837	0.00346501	6.200	8.46e-010	* * *				
exper	0.0214837	0.00346501	6.200	8.46e-010	* * *				

Standard error of the regression = 0.412948
#### 7.4 Testing a linear restriction

The calculated F-statistics (that used Equation 7.20) is 97.8892 with an associated p-value that is below 0.05. This means that the fit in the two regression equations is significantly different and we reject the null hypothesis presented in Equation 7.19. We conclude that the effect of education and experience have a significantly different effect on wages.

If you want to test whether education had four times the effect on wages than experience, the null is

$$H_0: \beta_2 = 4 \times \beta_3 \tag{7.21}$$

The command in Gretl is b[educ] - 4\*b[exper] = 0 to have

Standard error of the regression = 0.392836

Notice that the F-statistics is fairly small and has a p-value that is now greater than 5%. We do not reject the null hypothesis and conclude that, on average, one year of education has four times the effect on wages than one year of experience.

# Chapter 8 *Heteroscedasticity*

The fourth assumption in the estimation of the coefficients via ordinary least squares is the one of homoscedasticity. This means that the error terms  $u_i$  in the linear regression model have a constant variance across all observations *i*,

$$\sigma_{u_i}^2 = \sigma_u^2 \quad \text{for all } i. \tag{8.1}$$

When this assumption does not hold, and  $\sigma_{u_i}^2$  changes across *i* we say we have an heteroscedasticity problem. This chapter discusses the problems associated with heteroscedastic errors, presents some tests for heteroscedasticity and points out some possible solutions.

## 8.1 Heteroscedasticity and its implications

What happens if the errors are heteroscedasticity? The good news is that under heteroscedastic errors, OLS is still unbiased. The bad news is that we will obtain the incorrect standard errors of the coefficients. This means that the *t* and the *F* tests that we discussed in earlier chapters are no longer valid. Figure 8.1 shows the regression equation wage =  $\beta_0 + \beta_1$  educ + *u* with heteroscedastic errors. The variance of  $u_i$  increases with higher values of educ.

# 8.2 Testing for heteroscedasticity

#### 8.2.1 Breusch-Pagan test

Given the linear regression model

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_K + u$$
(8.2)

#### 8 Heteroscedasticity



**Fig. 8.1** wage =  $\beta_0 + \beta_1$  educ + *u* with heteroscedastic errors.

we know that OLS is unbiased and consistent if we assume  $E[u|X_2, X_3, ..., X_K] = 0$ . Let the null hypothesis that we have homoscedastic errors be

$$H_0: Var[u|X_2, X_3, \dots, X_K] = \sigma^2.$$
(8.3)

Because we are assuming that *u* has zero conditional expectation,  $Var[u|X_2, X_3, ..., X_K] = E[u^2|X_2, X_3, ..., X_K]$ , and so the null hypothesis of homoscedasticity is equivalent to

$$H_0: E[u^2|X_2, X_3, \dots, X_K] = \sigma^2.$$
(8.4)

This shows that if we want to test for violation of the homoscedasticity assumption, we want to test whether  $E[u^2|X_2, X_3, ..., X_K]$  is related to one or more of the independent variables. If  $H_0$  is false,  $E[u|X_2, X_3, ..., X_K]$  can be any function of the independent variables. A simple approach is to assume a linear function

$$u^2 = \delta_1 + \delta_2 X_1 + \delta_3 X_3 + \dots + \delta_K X_K + \varepsilon, \qquad (8.5)$$

where  $\varepsilon$  is an error term with mean zero given  $X_2, X_3, \ldots, X_K$ . The null hypothesis for homoscedasticity is:

$$H_0: \delta_1 = \delta_2 = \delta_3 = \dots = \delta_K = 0. \tag{8.6}$$

Under the null, it is reasonable to assume that  $\varepsilon$  is independent of  $X_2, X_3, \ldots, X_K$ . To be able to implement this test, we follow a two step procedure. In the first step we estimate Equation 5.14 via OLS. We estimate the residuals, square them and then

#### 8.2 Testing for heteroscedasticity

estimate the following equation:

$$\hat{u}^2 = \delta_1 + \delta_2 X_1 + \delta_3 X_3 + \dots + \delta_K X_K + error. \tag{8.7}$$

We can then easily compute the *F* statistic for the joint significance of all variables  $X_2, X_3, \ldots, X_K$ . Using OLS residuals in place of the errors does not affect the large sample distribution of the *F* statistic. An additional *LM* statistic to test for heteroscedasticity can be constructed based on the  $R_{a^2}^2$  obtained from Equation 8.7:

$$LM = n \cdot R_{\mu^2}^2. \tag{8.8}$$

Under the null hypothesis, *LM* is distributed asymptotically as  $\chi^2_{K-1}$ . This *LM* version of the test is called the Breusch-Pagan test for heteroscedasticity.

#### 8.2.2 Breusch-Pagan test in Gretl

As an example, consider once again our wage equation

wage = 
$$\beta_1 + \beta_2$$
educ +  $u$  (8.9)

Once we estimated the model in Gretl

T W

$$\widehat{wage} = \underbrace{146.952}_{(77.715)} + \underbrace{60.2143}_{(5.6950)} \text{ educ}$$

$$N = 935 \quad \overline{R}^2 = 0.1060 \quad F(1,933) = 111.79 \quad \widehat{\sigma} = 382.32$$
(standard errors in parentheses)

In the regression output window, go to  $\texttt{Tests} \to \texttt{Heteroskedasticity} \to \texttt{Breusch-Pagan}$  to obtain

```
Breusch-Pagan test for heteroskedasticity
OLS, using observations 1-935
Dependent variable: scaled uhat<sup>2</sup>
```

	coefficient	std. error	t-ratio	p-value	
const educ	-0.885844 0.140019	0.450097 0.0329833	-1.968 4.245	0.0494 2.40e-05	* * * * *
Explained	sum of squares	= 88.3581			
est statist ith p-value	ic: LM = $44.17$ = P(Chi-squar	9066, e(1) > 44.179	9066) = 0.0	00000	

Notice how Gretl reports the auxiliary regression presented in Equation 8.7 and the *LM* statistic from Equation 8.8. The large *LM* statistic associated with a small p-value (below 0.05 or 5%) indicates that we reject the null hypothesis of homoscedasticity. Hence, we have heteroscedaticity in the model of Equation 8.9.

#### 8.2.3 White test

White (1980) proposed a test for heteroscedasticity that that adds the squares and cross products of all the independent variables to Equation 8.2. In a model with only three independent variables, the White test is based on the estimation of:

$$\hat{u}^{2} = \delta_{1} + \delta_{2}X_{2} + \delta_{3}X_{3} + \delta_{4}X_{4} + \delta_{5}X_{2}^{2} + \delta_{6}X_{3}^{2} + \delta_{7}X_{4}^{2}$$

$$\delta_{8}X_{2} \cdot X_{3} + \delta_{9}X_{2} \cdot X_{4} + \delta_{10}X_{3} \cdot X_{4} + error.$$
(8.10)

Compared with the Breusch-Pagan test (see Equation 8.7), Equation 8.10 has more regressors. The White test for heteroscedasticity is based on the *LM* statistic for testing that all the  $\delta_i$  in Equation 8.10 are zero, except for the intercept.

## 8.2.4 White test in Gretl

We not use Gretl to test for heteroscedasticity in Equation 8.9 using the White test. In the regression output window, go to Tests  $\rightarrow$  Heteroskedasticity  $\rightarrow$  White's test to obtain

```
White's test for heteroskedasticity
OLS, using observations 1-935
Dependent variable: uhat^2
```

	coefficient	std. error	t-ratio	p-value
const educ	-126650 20049.7	435765 63065.6	-0.2906 0.3179	0.7714
sq_educ	13.2563	2234.63	0.005932	0.9953

Unadjusted R-squared = 0.018950

Test statistic:  $TR^2 = 17.717812$ , with p-value = P(Chi-square(2) > 17.717812) = 0.000142

Consistent with the Breusch-Pagan test, here the White test has a large *LM* statistic (labeled TR<sup>2</sup> following  $LM = n \cdot R_{u^2}^2$  as in Equation 8.8) associated with a small p-value (smaller than 5%). Hence, we reject the null of homoscedasticity and conclude that our model is heteroscedastic.

#### 8.3 What to do with heteroscedasticity?

There is a number of possible solutions when heteroscedastic errors are found. This section proposes three ways to solve the heteroscedasticity problem. First, a simple transformation of the variables; second, the use of weighted least squares; and third, the use of heteroscedasticity-robust standard errors.

8.3 What to do with heteroscedasticity?

#### 8.3.1 Simple transformation of the variables

An easy way to obtain homoscedastic errors is to come up with a simple transformation of the variables. Let's revisit the estimation of Equation 8.9, but this time with a simple logarithm transformation of wages,

$$\log wage = \beta_1 + \beta_2 educ + u \tag{8.11}$$

The Gretl regression output is

$$logwage = 5.97306 + 0.0598392 \text{ educ}$$

$$(0.081374) = (0.0059631)$$

$$N = 935 \quad \bar{R}^2 = 0.0964 \quad F(1,933) = 100.70 \quad \hat{\sigma} = 0.40032$$
(standard errors in parentheses)

Now, if we want to test for the existence of heteroscedasticity we go to  $\texttt{Tests} \rightarrow \texttt{Heteroskedasticity} \rightarrow \texttt{Breusch-Pagan}$  to obtain

Notice that the p-value associated with this test is above 0.05. Hence, we fail to reject the null of homoscedasticity. Compare this homoscedasticity results with the heteroscedastic errors found earlier when the variable wage was not in logs.

#### 8.3.2 Weighted Least Squares

We want to estimate the following regression model

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_K X_K + u, \tag{8.12}$$

but the errors u are heteroscedastic. When one is willing to assume that the heteroscedasticity appears as some function of  $X_2, X_3, \ldots, X_K$ , one can use Weighted Least Squares (WLS) to obtain homoscedastic errors. Let's say that the variance of u can be approximated using

8 Heteroscedasticity

$$u^2 = \sigma^2 \exp(\delta_1 + \delta_2 X_2 + \delta_3 X_3 + \dots + \delta_K X_K)\eta, \qquad (8.13)$$

where  $\eta$  is a random variable with mean equal to unity. If we assume that  $\eta$  is independent from  $X_2, X_3, \ldots, X_K$  we have

$$\log(u^2) = \alpha_1 + \delta_2 X_2 + \delta_3 X_3 + \dots + \delta_K X_K + \varepsilon.$$
(8.14)

To be able to implement this procedure, we replace the unobserved u with the OLS estimated residuals  $\hat{u}$  to estimate:

$$\log(\hat{u}^2) = \alpha_1 + \delta_2 X_2 + \delta_3 X_3 + \dots + \delta_K X_K + \varepsilon.$$
(8.15)

Finally, once Equation 8.15 is estimated, we obtain the fitted values and calculate the exponent to obtain

$$\hat{h}_i = \exp(\log(\hat{u}^2)). \tag{8.16}$$

We can use this  $\hat{h}_i$  as a weight in a Weighted Least Squares regression to solve the heteroscedasticity problem. That is, we estimate the following weighted equation

$$\frac{Y}{\hat{h}} = \beta_1 \frac{1}{\hat{h}} + \beta_2 \frac{X_2}{\hat{h}} + \beta_3 \frac{X_3}{\hat{h}} + \dots + \beta_K \frac{X_K}{\hat{h}} + \frac{u}{\hat{h}}.$$
(8.17)

Notice that Equation 8.17 is just Equation 8.12 divided by the weight  $\hat{h}_i$ . The new error term  $u/\hat{h}$  should be homoscedastic.

### 8.3.3 Weighted Least Squares in Gretl

Consider the following model

$$sav = \beta_1 + \beta_2 inc + u \tag{8.18}$$

where sav is savings and inc is income. The Gretl output is

Model 1: OLS, Dependent var	using of the state	observa sav	tions	s 1-10	00			
	coeffic	ient	std.	error	2	t-ratio	p-value	9
const inc	124.842 0.146	628	655.3 0.0	93) 57548	38	0.1905 2.548	0.8493 0.0124	**
Mean depender Sum squared m R-squared F(1, 98)	nt var resid	1582.5 1.00e+ 0.0621 6.4917	10 09 27 78	S.D. S.E. Adjus P-val	depe of r sted Lue(F	endent var regression R-squared ')	3284. 3197. 0.052 0.012	.902 .415 2557 2391
Log-likelihoo Schwarz crite	od · erion	947.89- 1904.9	935 97	Akaik Hanna	ke cr an-Qu	iterion inn	1899. 1901.	.787 .896

and the Breusch-Pagan test for heteroscedasticity yields

#### 8.3 What to do with heteroscedasticity?

```
Breusch-Pagan test for heteroskedasticity
OLS, using observations 1-100
Dependent variable: scaled uhat<sup>2</sup>
```

	coefficient	std. error	t-ratio	p-value
const	0.0457266	1.14381	0.03998	0.9682
inc	9.59914e-05	0.000100436	0.9557	0.3416

Explained sum of squares = 28.444

Test statistic: LM = 14.221987, with p-value = P(Chi-square(1) > 14.221987) = 0.000162

That is, we have heteroscedastic errors.

To estimate the WLS regression

$$\frac{\text{sav}}{\hat{h}} = \beta_1 \frac{1}{\hat{h}} + \beta_2 \frac{\text{inc}}{\hat{h}} + \frac{u}{\hat{h}}, \qquad (8.19)$$

in the Gretl main window we have to go to Model  $\rightarrow$  Other linear models  $\rightarrow$  Heteroskedasticity corrected to get the following computer output

Model 2: Heteroskedasticity-corrected, using observations 1-100 Dependent variable: sav

	coefficient	std. error	t-ratio	p-value	
const	-233.130	460.844	-0.5059	0.6141	
inc	0.185993	0.0616965	3.015	0.0033 **	*

Statistics based on the weighted data:

Sum squared resid	1043.864	S.E. of regression	3.263689
R-squared	0.084866	Adjusted R-squared	0.075527
F(1, 98)	9.088089	P-value(F)	0.003276
Log-likelihood	-259.1695	Akaike criterion	522.3391
Schwarz criterion	527.5494	Hannan-Quinn	524.4478

Statistics based on the original data:

Mean dependent var	1582.510	S.D. dependent var	3284.902
Sum squared resid	1.01e+09	S.E. of regression	3205.216

```
\widehat{\text{sav}} = -233.130 + 0.185993 \text{ inc} \\ (460.84) \quad (0.061697) \\ N = 100 \quad \overline{R}^2 = 0.0755 \quad F(1,98) = 9.0881 \quad \widehat{\sigma} = 3.2637 \\ (\text{standard errors in parentheses})
```

The estimates of the standard errors can now be used for inferences. The statistically significant coefficient on inc indicates that the marginal propensity to save out of

your income is 0.18. Of every additional dollar that you make, you will save 18 cents.

#### 8.3.4 White's heteroscedasticity-consistent standard errors

Even under the presence of heteroscedastic errors, at least in large samples a consistent estimator of the variances of the coefficients can be obtained via White's heteroscedasticity-consistent standard errors. This procedure leaves the OLS coefficients unaffected. For the estimation of Equation 8.18 you just have to make sure to select the option Robust standard errors in the Gretl "specify model" window when you estimate the model via OLS

Notice that the constant and slope coefficients are the same as before. However, the estimated standard errors are different.

# Chapter 9 Binary Choice Models

Some time we are interested in analyzing *binary response* or *qualitative response variables* that have outcomes *Y* equal to 1 when the even occurs and equal to 0 when the event does not occur. Some example include going to college, getting married, buying a house, or getting a job. All these cases involve a yes/no answer. How is this yes/no answer affected by other variables? That is the subject matter of this chapter.

## 9.1 The linear probability model

### 9.1.1 The model

The simplest binary choice model is the *linear probability model*, where as its name suggests, the probability of the event occurring, *p*, is assumed to be a linear function of a set of explanatory variable. If we only have one variable the model is

$$p_i = p(Y_i = 1) = \beta_1 + \beta_2 X_i.$$
 (9.1)

The response variable  $Y_i$  can be written as the summation of its deterministic and its random component,

$$Y_i = E(Y_i|X_i) + u_i.$$
 (9.2)

It is simple to compute  $E(Y_i|X_i)$ , the expected value of  $Y_i$  given  $X_i$ , because Y takes only two values. It is 1 with probability  $p_i$  and 0 with probability  $1 - p_i$ ,

$$E(Y_i|X_i) = 1 \times p_i + 0 \times (1 - p_i) = p_i = \beta_1 + \beta_2 X_i.$$
(9.3)

This means that we can write the model as

$$Y_i = \beta_1 + \beta_2 X_i + u_i, \tag{9.4}$$

that is just the same model we have been estimating in previous chapters. The big difference is that  $Y_i$  takes only the values 0 and 1.

### 9.1.2 The linear probability model in Gretl

Let's estimate the following model

$$\operatorname{inlf}_{i} = \beta_{1} + \beta_{2} \operatorname{educ}_{i} + \beta_{3} \operatorname{faminc}_{i} + u_{i}, \qquad (9.5)$$

where inlf is equal to one if individual *i* is in the labor force, zero otherwise, educ is the number of years of education and faminc is the family income. The regression command in Gretl is the same as before.

```
Model 1: OLS, using observations 1-753 Dependent variable: inlf
```

	coeffic	Lent	std.	error	t-ratio	p-value	
const educ faminc	0.068988 0.03790 1.459406	37 10 e-06	0.097 0.008 1.563	3736 35610 306e-06	0.7085 4.536 0.9337	0.4789 6.67e-06 → 0.3508	* * '
Mean depender Sum squared n R-squared F(2, 750) Log-likelihoo Schwarz crite	nt var cesid od - erion	0.5683 178.03 0.0362 14.093 -525.51 1070.9	393 367 221 349 .92 911	S.D. depe S.E. of m Adjusted P-value(E Akaike cm Hannan-Qu	endent var regression R-squared 7) riterion ainn	0.495630 0.487219 0.033651 9.81e-07 1057.038 1062.383	) ) L 7 3
$\widehat{inlf} = 0$	.0689887 (0.097374)	+0.037 (0.008	9040∈ 3561)	duc + 1.45	5940e-006 fa 5631e-006)	aminc	

N = 753  $\bar{R}^2 = 0.0337$  F(2,750) = 14.093  $\hat{\sigma} = 0.48722$ (standard errors in parentheses)

For example, the coefficient on educ indicates that every additional years of education increases the probability of being in the labor force by about 4%. This information is graphed in Figure 9.1.

There are two main problems with a linear probability model such as the one presented in Equation 9.4.

- 1. The model will predict unrealistic probabilities, beyond 1 and below 0 (see Figure 9.1).
- 2. Because  $Y_i$  only takes the values of 0 and 1, the error term u will be far from following a normal distribution.

The solution is to transform the linear probability model. Two common transformations are the *logit* and the *probit*.

#### 9.2 Logit analysis



**Fig. 9.1**  $inlf_i = \beta_1 + \beta_2 educ_i + \beta_3 faminc_i + u_i$ 

## 9.2 Logit analysis

# 9.2.1 The logit transformation

Let  $Z_i$  be,

$$Z_i = \beta_1 + \beta_2 X_i \tag{9.6}$$

The logit model hypothesizes that the probability of occurrence of the event Y = 1 is determined by the function

$$p_i = F(Z_i) = \frac{1}{1 + e^{-Z}} \tag{9.7}$$

where

$$\frac{\partial p}{\partial Z} = f(Z) = \frac{e^{-Z}}{(1+e^{-Z})^2}$$
(9.8)

and

$$\frac{\partial p}{\partial X} = \frac{\partial p}{\partial Z} \frac{\partial Z}{\partial X} = f(Z)\beta_2 \tag{9.9}$$

This means that the marginal effect of variable X on the probability of Y = 1 is  $f(Z)\beta_2$ , where f(Z) needs to be evaluated on some specific value of X, let's say, the mean of X.

#### 9.2.2 Logit regression in Gretl

Fortunately, all these calculations are done automatically by Gretl. If we want to obtain the logit estimates of Equation 9.5 in the main Gretl window we have to go to Model  $\rightarrow$  Nonlinear models  $\rightarrow$  Logit  $\rightarrow$  Binary... and select the option "Show p-values" to obtain

Convergence achieved after 4 iterations Model 2: Logit, using observations 1-753 Dependent variable: inlf coefficient std. error z p-value \_\_\_\_\_ const-1.852870.428444-4.3251.53e-05\*\*\*educ0.1617730.03678564.3981.09e-05\*\*\*faminc6.58050e-066.83134e-060.96330.3354 Mean dependent var 0.568393 S.D. dependent var 0.244933 McFadden R-squared 0.027185 Adjusted R-squared 0.021359 Log-likelihood -500.8762 Akaike criterion 1007.752 Schwarz criterion 1021.625 Hannan-Quinn 1013.097 Number of cases 'correctly predicted' = 449 (59.6%) f(beta'x) at mean of independent vars = 0.245Likelihood ratio test: Chi-square(2) = 27.9939 [0.0000] Predicted A

	0	1
0	69	256
1	48	380
	0 1	0 0 1 48

Gretl actually estimates this model using an estimation technique called Maximum Likelihood Estimation, that is why the computer iterates before giving the estimates. The output is very similar as the one obtained in previous chapters. The effect of educ on inlf is statistically significant. However, the key difference in this output is that the coefficients are not interpreted as the marginal effects. Recall that the marginal effects are calculated using Equation 9.9. To make Gretl obtain this marginal effects you need to reestimate the model and select the option "Show slopes at mean" to obtain

	coefficient	std. error	Z	slope
const	-1.85287	0.428444	-4.325	0 0396234
faminc	6.58050e-06	6.83134e-06	0.9633	1.61178e-06

The marginal effect of educ on inlf is actually 0.0396. An additional year of education will increase the probability that you are in the labor force by about 4%.

80

9.3 Probit analysis

#### **9.3 Probit analysis**

#### 9.3.1 The probit transformation

The probit model is similar in spirit as the logit model. Let  $Z_i$  be,

$$Z_i = \beta_1 + \beta_2 X_i \tag{9.10}$$

The probit model hypothesizes that the probability of occurrence of the event Y = 1 is determined by the function

$$p_i = F(Z_i) \tag{9.11}$$

where  $F(\cdot)$  is actually the cumulative standardized normal distribution. Then,

$$\frac{\partial p}{\partial Z} = f(Z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}Z^2}$$
(9.12)

is just the derivative of  $F(\cdot)$ . As in the logit case,

$$\frac{\partial p}{\partial X} = \frac{\partial p}{\partial Z} \frac{\partial Z}{\partial X} = f(Z)\beta_2 \tag{9.13}$$

Again, this means that the marginal effect of variable X on the probability of Y = 1 is  $f(Z)\beta_2$ , where f(Z) needs to be evaluated on some specific value of X, let's say, the mean of X.

#### 9.3.2 Probit regression in Gretl

If we want to obtain the probit estimates of Equation 9.5 in the main Gretl window we have to go to Model  $\rightarrow$  Nonlinear models  $\rightarrow$  Probit  $\rightarrow$  Binary... and select the option "Show p-values" to obtain

```
Number of cases 'correctly predicted' = 449 (59.6%)
f(beta'x) at mean of independent vars = 0.393
Likelihood ratio test: Chi-square(2) = 28.0252 [0.0000]

Predicted
0 1
Actual 0 69 256
1 48 380
```

Once again, the effect of educ on inlf is statistically significant. To make Gretl obtain this marginal effects using Equation 9.13 you need to reestimate the model and select the option "Show slopes at mean" to obtain

	coefficient	std. error	Z	slope
const educ faminc	-1.14768 0.100666 3.84752e-06	0.261470 0.0224296 4.09647e-06	-4.389 4.488 0.9392	0.0395289 1.51081e-06

The marginal effect of educ on inlf is 0.0395. We obtain almost the same results as before. The probit model predicts that an additional year of education will increase the probability that you are in the labor force by about 4%.

82

# Chapter 10 *Time Series*

## **10.1 Time Series Data**

The main difference between time series data and cross-sectional data is the temporal ordering. To emphasize the proper ordering of the observations, Table 10.1 presents a partial listing of the data on U.S. inflation and unemployment rates from 1948 through 2003. Unlike cross-sectional data, in time series the temporal order in which the observations appear in the data set is very important. In terms of notation, we use the subscript *t* to denote time and we use it instead of the previous subscript *i*, i.e., *X<sub>t</sub>*.

Year	Inflation	Unemployment
1948	8.1	3.8
1949	-1.2	5.9
1950	1.3	5.3
1951	7.9	3.3
:	•	:
2000	3.4	4.0
2001	2.8	4.7
2002	1.6	5.8
2003	2.3	6.0

Table 10.1 U.S. Inflation and Unemployment Rates, 1965-2011

A second key difference between time series and cross-sectional data is that in the latter we assume that the sample was randomly drawn from the population. While in time series the variables are also considered random, a variable indexed by time is called a *stochastic process* or a *time series process*. When we collect a time series data set we are one possible outcome or realization of the stochastic process. We can only see a single realization because we cannot go back in time and start the process again.



Fig. 10.1 Inflation, 1948-2003.

Graphing the data is particularly important to visualize the dynamics of the variables. Figure 10.1 presents the time series graph of inflation from 1948 through 2003. One can easily identify the periods of high inflation late in the seventies and early eighties. To obtain this graph in Gretl, go to  $View \rightarrow Graph$  specified vars  $\rightarrow$  Time series plot and then select the variables you want to plot against time.

## **10.2 Time Series Regression Models**

# 10.2.1 Static Models

The simplest static model has the form

$$Y_t = \beta_1 + \beta_2 X_t + u_t, \quad t = 1, 2, 3, \dots, n.$$
(10.1)

#### 10.2 Time Series Regression Models

We call this a static model because we are only modeling a contemporaneous relationship between  $X_t$  and  $Y_t$ . That is, when a change in X at time t is believed to have an immediate effect on  $Y: \Delta Y_t = \beta_2 \Delta X_t$ . One example is the static Phillips curve given by:

inflation<sub>t</sub> = 
$$\beta_1 + \beta_2$$
unemployment<sub>t</sub> +  $u_t$ . (10.2)

where inflation is the annual inflation rate, and unemployment is the unemployment rate. Estimation in Gretl via OLS is follows the same steps as in the previous chapters. The output for the estimation of Equation 10.1 is:

```
Model 1: OLS, using observations 1948-2003 (T = 56) Dependent variable: inflation
```

	coefficient	std. error	t-ratio	p-value
const	1.05357	1.54796	0.6806	0.4990
unemployment	0.502378	0.265562	1.892	0.0639 *
Mean dependent	var 3.8839	29 S.D. dep	pendent var	3.040381
Sum squared re	sid 476.81	57 S.E. of	regression	2.971518
R-squared	0.0621	54 Adjusted	d R-squared	0.044786
F(1, 54)	3.5787	26 P-value	(F)	0.063892
Log-likelihood	-139.43	04 Akaike o	criterion	282.8607
Schwarz criter	ion 286.91	14 Hannan-(	Quinn	284.4311
rho	0.5720	55 Durbin-W	Natson	0.801482

 $\widehat{\text{inflation}} = 1.05357 + 0.502378 \, \text{unemployment}_{(1.5480)} + 0.26556)$ 

$$T = 56$$
  $\bar{R}^2 = 0.0448$   $F(1,54) = 3.5787$   $\hat{\sigma} = 2.9715$ 

(standard errors in parentheses)

The estimation results indicate that a one point increase in the unemployment rate is linked with a 0.5 increase in the inflation rate. Of course more variables can be included in the model. Notice that we can use this model to predict inflation given that we know the values for unemployment by simply plugging values for unemployment in the estimated equation. If we do this for the actual unemployment values for 1948-2003 period and graph them, we obtain the fitted values graph. Figure 10.2 plots the actual and the fitted values for inflation.

# 10.2.2 Finite Distributed Lag Models

The simplest dynamic model is the *finite distributed lag* (FDL) model, where we allow one or more variables to to affect  $Y_t$  with a lag. Consider the following example:

$$Y_t = \beta_1 + \beta_2 X_t + \beta_3 X_{t-1} + \beta_4 X_{t-2} + u_t, \qquad (10.3)$$



Fig. 10.2 Inflation, 1948-2003. Actual and fitted based on an static model.

where the FDL is of order two. Let's say that we are interested in the effect on *Y* of a permanent increase in *X*. Before time t, *X* equals to a constant *c*. At time *t*, *X* increases permanently to c + 1. That is,  $X_s = c$  for s < t and  $X_s = c + 1$  for  $s \ge t$ . Setting the errors to be equal to zero we have:

$$Y_{t-1} = \beta_1 + \beta_2 c + \beta_3 c + \beta_4 c$$

$$Y_t = \beta_1 + \beta_2 (c+1) + \beta_3 c + \beta_4 c$$

$$Y_{t+1} = \beta_1 + \beta_2 (c+1) + \beta_3 (c+1) + \beta_4 c$$

$$Y_{t+2} = \beta_1 + \beta_2 (c+1) + \beta_3 (c+1) + \beta_4 (c+1)$$
(10.4)

and so on. The contemporaneous effect of X on Y is called the *impact multiplier* and in this case this one is given by  $\beta_2$ . However, over time the marginal effect of X on Y is larger. We say that the *long-run multiplier* is the long-run change Y given a permanent increase in X. This one is given by  $\beta_2 + \beta_3 + \beta_4$ .

Consider the following example in Gretl:

$$\inf_t = \beta_1 + \beta_2 \operatorname{unem}_t + \beta_3 \operatorname{unem}_{t-1} + \beta_4 \operatorname{unem}_{t-2} + u_t, \quad (10.5)$$

10.2 Time Series Regression Models

To estimate this model in Gretl we first need to create the lagged values of unem. To do this we have to go to select unem and then go to Add  $\rightarrow$  Lags of selected variables and select the number of lags. An alternative approach is to just include the lags when estimating the model via OLS. That is, when specifying the model in Gretl (Model  $\rightarrow$  Ordinary Least Squares) there is an icon that allows you to select the lags. Just select two lags for unem to obtain:

Model 2: OLS, using observations 1950-2003 (T = 54) Dependent variable: inf

	coeffic	cient	std.	error	t-ratio	p-value	9
const	-0.124	1609	1.6	8922	-0.07377	0.9415	-
unem	0.903	3211	0.4	02071	2.246	0.0291	* *
unem_1	-0.850	5337	0.5	25700	-1.629	0.1096	
unem_2	0.668	8123	0.3	86722	1.728	0.0902	*
Mean depende	nt var	3.900	0000	S.D. (	dependent var	2.961	.323
Sum squared	resid	395.2	2340	S.E.	of regression	2.811	526
R-squared		0.149	9632	Adjust	ted R-squared	0.098	8610
F(3, 50)		2.932	2693	P-val	ue(F)	0.042	366
Log-likeliho	od	-130.3	3660	Akaik	e criterion	268.7	320
Schwarz crit	erion	276.0	6880	Hannai	n-Quinn	271.8	8003
rho		0.661	1217	Durbi	n-Watson	0.676	5987

$$\widehat{\ln f} = -0.124609 + 0.903211 \text{ unem} - 0.856337 \text{ unem} + 0.668123 \text{ unem} - 2.8570}$$
$$T = 54 \quad \overline{R}^2 = 0.0986 \quad F(3,50) = 2.9327 \quad \widehat{\sigma} = 2.8115$$
(standard errors in parentheses)

A permanent increase in unemployment leads to a contemporaneous increase in inflation of 0.903 (impact multiplier). However, in the long-run the same increase in unemployment leads to a permanent effect on inflation of 0.903 - 0.856 + 0.668 = 0.715 (long-run multiplier).

## **10.2.3** Autoregressive Model

An *autoregresive model* is a simple model where the current values of a variable are related to its past values. The first-order autoregressive model is given by:

$$Y_t = \phi Y_{t-1} + u_t. \tag{10.6}$$

This one is usually denoted by AR(1). A more general model is the *p*th autoregressive model or AR(p) given by:

$$Y_t = \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \phi_3 Y_{t-3} + \dots + \phi_p Y_{t-p} + u_t$$
(10.7)

where there are p lags of the variable Y explaining its current value. The estimation of an AR(p) model in Gretl is simple; go to Model  $\rightarrow$  Time series  $\rightarrow$  ARIMA and then select the dependent variable and the AR order. Make sure that the MAorder is zero. For the example above, consider estimating the following model:

$$\inf_{t} = \phi_1 \inf_{t-1} + \phi_2 \inf_{t-2} + u_t \tag{10.8}$$

The output in Gretl is:

Mean depe	ndent var	3.883929	S.D. depe	endent var	3.040381
Mean of i	nnovations	-0.054104	S.D. of i	nnovations	2.171763
Log-likelihood		-123.2506	Akaike cr	254.5012	
Schwarz c	riterion	262.6026	Hannan-Qu	linn	257.6421
		Real	Imaginary	Modulus	Frequency
AR					
Root	1	2.3276	-0.5378	2.3889	-0.0361
Root	2	2.3276	0.5378	2.3889	0.0361

That yields the following estimated equation:

$$\inf_{t} = 4.025 + 0.8157 \inf_{t-1} - 0.1752 \inf_{t-2}.$$
 (10.9)

where we can see that higher inflation last period has a positive effect on inflation this period. We can use this model to predict the path of inf based on its previous values. It Gretl the command to obtain this Graphs  $\rightarrow$  Fitted, actual plot  $\rightarrow$  Against time. The resulting graph is shown in Figure 10.3.

### 10.2.4 Moving-Average Models

The *moving-average* models express an observed series as a function of the current and lagged unobserved shocks. The simplest moving-average model is the moving-average of order one, or MA(1):



Fig. 10.3 Inflation, 1948-2003. Actual and fitted based on an AR(2) model.

$$Y_t = \theta u_{t-1} + u_t \tag{10.10}$$

A more general moving-average of order q is be written as:

$$Y_t = \theta_1 u_{t-1} + \theta_2 u_{t-2} + \theta_3 u_{t-3} + \dots + \theta_q u_{t-q} + u_t$$
(10.11)

For the example above:

$$\inf_{t} = \theta_{1} u_{t-1} + \theta_{2} u_{t-2} + u_{t} \tag{10.12}$$

the output in Gretl is:

```
Function evaluations: 51
Evaluations of gradient: 19
Model 6: ARMA, using observations 1948-2003 (T = 56)
Estimated using Kalman filter (exact ML)
```

```
Dependent variable: inf
Standard errors based on Hessian
```

coefficient std. error z p-value



Fig. 10.4 Inflation, 1948-2003. Actual and fitted based on an MA(2) model.

const	3.9820	57 0	.615240	6.473	9.58	8e-011	***
theta_1	1.1854	19 0	.130300	9.098	9.19	e-020	***
theta_2	0.2679	922 0	.127516	2.101	0.03	356	**
Mean depe	ndent var	3.883929	S.D.	dependent	var	3.040	)381
Mean of i	nnovations	-0.041523	S.D.	of innovat	ions	1.899	9580
Log-likel	ihood	-116.5030	Akaik	ke criteric	n	241.0	061
Schwarz c	riterion	249.1075	Hanna	an-Quinn		244.1	L470
		Real	Imagina	ary Modu	lus	Freque	ency
MA							
Root	1	-1.1343	0.00	000 1.1	343	0.5	5000
Root	2	-3.2904	0.00	)00 3.2	904	0.5	5000

and the actual and fitted values are presented in Figure 10.4.

10.2 Time Series Regression Models

## 10.2.5 Autoregressive Moving Average Models

One can easily combine an AR(1) model and an MA(1) models to obtain an autoregressive moving-average model ARMA(1,1):

$$Y_t = \phi Y_{t-1} + \theta u_{t-1} + u_t \tag{10.13}$$

or a more general ARMA(p,q) model:

$$Y_t = \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \theta_1 u_{t-1} + \theta_2 u_{t-2} + \dots + \theta_q u_{t-q} + u_t \quad (10.14)$$

The output in Gretl for a ARMA(2,2) for inflation is:

```
Function evaluations: 61
Evaluations of gradient: 20
Model 5: ARMA, using observations 1948-2003 (T = 56)
Estimated using Kalman filter (exact ML)
Dependent variable: inf
Standard errors based on Hessian
```

	coefficient	std. erro	r z	p-value	è
const phi_1 phi_2 theta_1 theta_2	3.94843 0.828806 0.0226838 0.274108 -0.587919	1.05291 0.236639 0.173277 0.197397 0.169467	3.750 3.502 0.1309 1.389 -3.469	0.0002 0.0005 0.8958 0.1650 0.0005	- *** ***
cilecta_2	0.307919	0.109407	5.105	0.0005	~ ~ ~ ~
Mean depender Mean of innov Log-likelihoo Schwarz crite	nt var 3.883 vations -0.053 od -114.4 erion 253.1	929 S.D. 524 S.D. 824 Akai 169 Hann	dependent of innovat ke criterio an-Quinn	var 3.04 ions 1.83 n 240. 245.	0381 1760 9648 6761
	Re	al Imagin	ary Modu	lus Frequ	lency
 AR					
Root 1	1.16	91 0.0	000 1.1	691 0.	0000
Root 2	-37.70	65 0.0	000 37.7	065 0.	5000
MA					
Root 1	-1.09	17 0.0	000 1.0	917 0.	5000
Root 2	1.55	80 0.0	000 1.5	580 0.	0000