

# Lecture 2: Simple Linear Regression

R.G. Pierse

## 1 The Classical Linear Regression Model

$$Y_i = a + bX_i + u_i \quad , \quad i = 1, \dots, n \quad (1.1)$$

### 1.1 Assumptions of the Model

$$E(u_i) = 0 \quad , \quad i = 1, \dots, n \quad (A1)$$

$$E(u_i^2) = \sigma^2 \quad , \quad i = 1, \dots, n \quad (A2)$$

$$E(u_i u_j) = 0 \quad , \quad i, j = 1, \dots, n \quad j \neq i \quad (A3)$$

$$X \text{ values are } \textit{fixed in repeated sampling} \quad (A4)$$

$Y_i$  is a random variable with the following properties:

$$E(Y_i) = a + bX_i \quad , \quad i = 1, \dots, n$$

$$\begin{aligned} \text{Var}(Y_i) &= E(Y_i - E(Y_i))^2 \\ &= E(u_i^2) = \sigma^2 \quad , \quad i = 1, \dots, n \end{aligned}$$

## 2 The Ordinary Least Squares Estimator

$$Y_i = \hat{a} + \hat{b}X_i + e_i \quad , \quad i = 1, \dots, n \quad (2.1)$$

The OLS estimator is the estimator that minimises the sum of squared residuals  $s = \sum_{i=1}^n e_i^2$ .

$$\min_{\hat{a}, \hat{b}} s = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \hat{a} - \hat{b}X_i)^2$$

$$s = \sum_{i=1}^n Y_i^2 - 2\hat{a} \sum_{i=1}^n Y_i - 2\hat{b} \sum_{i=1}^n X_i Y_i + \sum_{i=1}^n \hat{a}^2 + 2\hat{a}\hat{b} \sum_{i=1}^n X_i + \hat{b}^2 \sum_{i=1}^n X_i^2$$

Differentiating this expression with respect to  $\hat{a}$  and  $\hat{b}$  gives the two first order conditions:

$$\frac{\partial s}{\partial \hat{a}} = -2 \sum_{i=1}^n Y_i + 2n\hat{a} + 2\hat{b} \sum_{i=1}^n X_i = 0 \quad (2.2)$$

and

$$\frac{\partial s}{\partial \hat{b}} = -2 \sum_{i=1}^n X_i Y_i + 2\hat{a} \sum_{i=1}^n X_i + 2\hat{b} \sum_{i=1}^n X_i^2 = 0 \quad (2.3)$$

Rearranging (2.2) gives

$$\hat{a} = \frac{1}{n} \sum_{i=1}^n Y_i - \hat{b} \frac{1}{n} \sum_{i=1}^n X_i = \bar{Y} - \hat{b} \bar{X} \quad (2.4)$$

where  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$  and  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  are the sample means of  $Y$  and  $X$  respectively.

Substituting this expression into (2.3) and cancelling the factor of 2 gives

$$\sum_{i=1}^n X_i Y_i - n\bar{Y} \bar{X} - \hat{b} n\bar{X}^2 - \hat{b} \sum_{i=1}^n X_i^2 = 0$$

which can be rearranged as:

$$\hat{b} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{Y} \bar{X}}{\sum_{i=1}^n X_i^2 - n\bar{X}^2}$$

or

$$\hat{b} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}. \quad (2.5)$$

Alternatively, using lower case characters to denote deviations from sample means, we can write

$$\hat{b} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}. \quad (2.6)$$

### 3 Properties of the OLS Estimator

Summing (1.1) and dividing by  $n$  gives

$$\bar{Y} = a + b\bar{X} + \bar{u}$$

and, subtracting from (1.1),

$$Y_i - \bar{Y} = b(X_i - \bar{X}) + (u_i - \bar{u})$$

or

$$y_i = bx_i + u_i - \bar{u}. \quad (3.1)$$

### 3.1 The OLS Estimators $\hat{a}$ and $\hat{b}$ are Unbiased

$$E(\hat{b}) = b \quad , \quad E(\hat{a}) = a$$

**Proposition 3.1.**  $E(\hat{b}) = b$

*Proof.* Substituting (3.1) into (2.6)

$$\hat{b} = b + \frac{\sum_{i=1}^n x_i(u_i - \bar{u})}{\sum_{i=1}^n x_i^2} \quad (3.2)$$

$$= b + \frac{\sum_{i=1}^n x_i u_i}{\sum_{i=1}^n x_i^2} \quad (3.3)$$

where the second equality follows because  $\sum_{i=1}^n x_i = 0$ .  
Taking expectations,

$$E(\hat{b}) = b + \frac{\sum_{i=1}^n x_i E(u_i)}{\sum_{i=1}^n x_i^2}$$

but, from assumption (A1), the second term is zero so that

$$E(\hat{b}) = b \quad (3.4)$$

and the estimator  $\hat{b}$  is *unbiased*. □

**Proposition 3.2.**  $E(\hat{a}) = a$

*Proof.* Substituting (1.1) into (2.4)

$$\begin{aligned} \hat{a} &= \frac{1}{n} \sum_{i=1}^n (a + bX_i + u_i) - \hat{b} \frac{1}{n} \sum_{i=1}^n X_i \\ &= a + (b - \hat{b}) \frac{1}{n} \sum_{i=1}^n X_i + \frac{1}{n} \sum_{i=1}^n u_i \end{aligned} \quad (3.5)$$

and, taking expectations,

$$E(\widehat{a}) = a + (b - E(\widehat{b})) \frac{1}{n} \sum_{i=1}^n X_i + \frac{1}{n} \sum_{i=1}^n E(u_i)$$

and the last two terms are zero, by (3.4) and (A1), so that

$$E(\widehat{a}) = a$$

and the estimator  $\widehat{a}$  is *unbiased*. □

### 3.2 The Variance of OLS Estimators

From (3.3)

$$\widehat{b} - E(\widehat{b}) = \frac{\sum_{i=1}^n x_i u_i}{\sum_{i=1}^n x_i^2}$$

so that

$$\text{Var}(\widehat{b}) \equiv E(\widehat{b} - E(\widehat{b}))^2 = E \left( \frac{\sum_{i=1}^n x_i u_i}{\sum_{i=1}^n x_i^2} \right)^2.$$

Expanding the numerator gives

$$\text{Var}(\widehat{b}) = \frac{\sum_{i=1}^n x_i^2 E(u_i^2) + \sum_{j \neq i} \sum_{i=1}^n x_i x_j E(u_i u_j)}{(\sum_{i=1}^n x_i^2)^2}$$

and from (A2) and (A3) it follows that

$$\text{Var}(\widehat{b}) = \frac{\sigma^2 \sum_{i=1}^n x_i^2}{(\sum_{i=1}^n x_i^2)^2} = \frac{\sigma^2}{\sum_{i=1}^n x_i^2} \tag{3.6}$$

Similarly, for  $\widehat{a}$ , from (3.5)

$$\widehat{a} - E(\widehat{a}) = (b - \widehat{b}) \frac{1}{n} \sum_{i=1}^n X_i + \frac{1}{n} \sum_{i=1}^n u_i$$

so that

$$\text{Var}(\widehat{a}) = E \left( (b - \widehat{b}) \frac{1}{n} \sum_{i=1}^n X_i + \frac{1}{n} \sum_{i=1}^n u_i \right)^2.$$

It can be proved (although here the result will only be stated) that

$$\text{Var}(\widehat{a}) = \sigma^2 \left( \frac{1}{n} + \frac{\overline{X}^2}{\sum_{i=1}^n x_i^2} \right) = \frac{\sigma^2 \sum_{i=1}^n X_i^2}{n \sum_{i=1}^n x_i^2} \tag{3.7}$$

### 3.3 The OLS estimator of $\sigma^2$

The error variance  $\sigma^2$  that appears in formulae (3.7) and (3.6) is itself unknown and so in practice we need to estimate it. We now show that the estimator

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n-2} \quad (3.8)$$

is an unbiased estimator of  $\sigma^2$ .

$$\begin{aligned} \sum_{i=1}^n e_i^2 &= \sum_{i=1}^n (Y_i - \hat{a} - \hat{b}X_i)^2 \\ &= \sum_{i=1}^n (y_i - \hat{b}x_i)^2 \end{aligned}$$

from (2.4). Substituting from (3.1)

$$\begin{aligned} \sum_{i=1}^n e_i^2 &= \sum_{i=1}^n (bx_i + u_i - \bar{u} - \hat{b}x_i)^2 \\ &= (\hat{b} - b)^2 \sum_{i=1}^n x_i^2 + \sum_{i=1}^n (u_i - \bar{u})^2 - 2(\hat{b} - b) \sum_{i=1}^n x_i(u_i - \bar{u}). \quad (3.9) \end{aligned}$$

Consider the last term in equation (3.9). Note that  $\sum_{i=1}^n x_i(u_i - \bar{u}) = \sum_{i=1}^n x_i u_i$  since  $\bar{u} \sum_{i=1}^n x_i = 0$  and, from equation (3.3)  $\sum_{i=1}^n x_i u_i = \sum_{i=1}^n x_i^2 (\hat{b} - b)$ , so that, substituting and taking expectations,

$$E\left(\sum_{i=1}^n e_i^2\right) = E(\hat{b} - b)^2 \sum_{i=1}^n x_i^2 + \sum_{i=1}^n E(u_i - \bar{u})^2 - 2E(\hat{b} - b)^2 \sum_{i=1}^n x_i^2. \quad (3.10)$$

From the expression for  $Var(\hat{b})$  (3.6) the first and third terms in (3.10) are  $\sigma^2$  and  $-2\sigma^2$  respectively. Expanding the second term

$$\begin{aligned} \sum_{i=1}^n E(u_i - \bar{u})^2 &= \sum_{i=1}^n E(u_i)^2 - 2E \sum_{i=1}^n (u_i \bar{u}) + nE(\bar{u}^2) \\ &= \sum_{i=1}^n E(u_i)^2 - 2nE(\bar{u}) + nE(\bar{u}^2) \\ &= \sum_{i=1}^n E(u_i)^2 - nE(\bar{u}^2) \quad (3.11) \end{aligned}$$

but

$$E(\bar{u}^2) = E\left(\frac{1}{n} \sum_{i=1}^n u_i\right)^2 = \frac{1}{n^2} E\left(\sum_{i=1}^n u_i^2 + \sum_{j \neq i} \sum_{i=1}^n u_i u_j\right) = \frac{1}{n} \sigma^2$$

so that, substituting into (3.11),

$$\sum_{i=1}^n E(u_i - \bar{u})^2 = \sum_{i=1}^n E(u_i^2) - nE(\bar{u}^2) = n\sigma^2 - \sigma^2 = (n-1)\sigma^2.$$

Finally, substituting these expressions back into (3.10), we have

$$E\left(\sum_{i=1}^n e_i^2\right) = \sigma^2 + (n-1)\sigma^2 - 2\sigma^2 = (n-2)\sigma^2$$

so that

$$E(\hat{\sigma}^2) = \frac{E\left(\sum_{i=1}^n e_i^2\right)}{(n-2)} = \sigma^2$$

which shows that  $\hat{\sigma}^2$  is an unbiased estimator of the error variance  $\sigma^2$ .

## 4 The Gauss-Markov Theorem

**Definition 4.1.** A linear estimator is one that can be written in the form

$$\tilde{b} = \sum_i^n w_i Y_i$$

where  $w_i$  are fixed weights.

Note that the OLS estimator  $\hat{b}$  is a linear estimator since

$$\hat{b} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} = \frac{\sum_{i=1}^n x_i Y_i}{\sum_{i=1}^n x_i^2} = \sum_{i=1}^n k_i Y_i$$

where the weights  $k_i$  are given by

$$k_i = \frac{x_i}{\sum_{i=1}^n x_i^2}.$$

**Theorem 4.1.** The OLS estimator  $\hat{b}$  is the Best Linear Unbiased Estimator (BLUE) of the classical regression model. By **best** we mean the estimator in the class that achieves **minimum variance**.

*Proof.* Taking expectations

$$\begin{aligned} E(\tilde{b}) &= \sum_i^n w_i E(Y_i) \\ &= a \sum_i^n w_i + b \sum_i^n w_i X_i \end{aligned}$$

so the conditions for unbiasedness of  $\tilde{b}$  are that

$$\sum_i^n w_i = 0 \quad \text{and} \quad \sum_i^n w_i X_i = 1 \quad (4.1)$$

The variance of the estimator  $\tilde{b}$  is given by

$$\text{Var}(\tilde{b}) = \sum_i^n w_i^2 \text{Var}(Y_i) = \sigma^2 \sum_i^n w_i^2$$

We now use a trick and add and subtract the OLS weights  $k_i$  from this expression to give

$$\begin{aligned} \text{Var}(\tilde{b}) &= \sigma^2 \sum_i^n (w_i - k_i + k_i)^2 \\ &= \sigma^2 \sum_i^n (w_i - k_i)^2 + \sigma^2 \sum_i^n k_i^2 + 2\sigma^2 \sum_i^n (w_i - k_i)k_i \end{aligned} \quad (4.2)$$

but the third term in (4.2) is zero from the unbiasedness conditions (4.1) since

$$\sum_i^n (w_i - k_i)k_i = \frac{\sum_{i=1}^n w_i X_i - \bar{X} \sum_{i=1}^n w_i}{\sum_{i=1}^n x_i^2} - \frac{1}{\sum_{i=1}^n x_i^2} = 0$$

Thus

$$\begin{aligned} \text{Var}(\tilde{b}) &= \sigma^2 \sum_i^n (w_i - k_i)^2 + \sigma^2 \sum_i^n k_i^2 \\ &= \sigma^2 \sum_i^n (w_i - k_i)^2 + \frac{\sigma^2}{\sum_i^n x_i^2} \\ &= \sigma^2 \sum_i^n (w_i - k_i)^2 + \text{Var}(\hat{b}) \end{aligned} \quad (4.3)$$

and the first term in (4.3) is greater than or equal to zero, achieving its lower bound for the OLS estimator where

$$w_i = k_i \quad \text{or equivalently} \quad \tilde{b} = \hat{b}.$$

□

## 5 Hypothesis Testing

In order to make statistical inferences on the parameter estimates  $\hat{a}$ ,  $\hat{b}$  and  $\hat{\sigma}^2$  we must make a further assumption:

$$u_i \sim iid N(0, \sigma^2) \quad , \quad i = 1, \dots, n \quad (\text{A5})$$

that the errors  $u_i$  are distributed as independent normal variables. It then follows that, since  $y_i$ ,  $\hat{a}$  and  $\hat{b}$  are all linear combinations of  $u_i$ , they are also distributed normally with

$$y_i \sim N(a + bX_i, \sigma^2) \quad , \quad i = 1, \dots, n$$

$$\hat{a} \sim N\left(a, \frac{\sigma^2 \sum_{i=1}^n X_i^2}{n \sum_{i=1}^n x_i^2}\right)$$

and

$$\hat{b} \sim N\left(b, \frac{\sigma^2}{\sum_{i=1}^n x_i^2}\right).$$

To make this practical we need to replace the unknown parameter  $\sigma^2$  with the estimator  $\hat{\sigma}^2$  defined in (3.8). Then, defining

$$\hat{\sigma}_{\hat{a}} = \sqrt{\frac{\hat{\sigma}^2 \sum_{i=1}^n X_i^2}{n \sum_{i=1}^n x_i^2}} \quad (5.1)$$

and

$$\hat{\sigma}_{\hat{b}} = \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n x_i^2}} \quad (5.2)$$

it can be shown that

$$\frac{\hat{a} - a}{\hat{\sigma}_{\hat{a}}} \sim t_{n-2} \quad (5.3)$$

and

$$\frac{\hat{b} - b}{\hat{\sigma}_{\hat{b}}} \sim t_{n-2} \quad (5.4)$$

where  $t_{n-2}$  is the Student t distribution with  $n - 2$  degrees of freedom.  $\hat{\sigma}_{\hat{a}}$  and  $\hat{\sigma}_{\hat{b}}$  are known as the *standard errors* of the estimators  $\hat{a}$  and  $\hat{b}$  respectively.

It can also be shown that

$$(n - 2) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2 \quad (5.5)$$

where  $\chi_{n-2}^2$  is the Chi-squared distribution with  $n - 2$  degrees of freedom.  $\hat{\sigma}$  is known as the *equation standard error*.