

## Ordinary Least Squares Regression.

### Simple Regression. Algebra and Assumptions.

In this part of the course we are going to study a technique for analysing the linear relationship between two variables Y and X. We have n pairs of observations  $(Y_i, X_i)$ ,  $i = 1, 2, \dots, n$  on the relationship which, because it is not exact, we shall write as:

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

In this relationship  $\alpha$ ,  $\beta$  and  $\epsilon_i$   $i = 1, \dots, n$  are fundamentally unobservable and we would like to estimate  $\alpha$  and  $\beta$ .

#### Approach:

The idea is to select estimates of  $\alpha$  and  $\beta$  lets call them  $\alpha^*$  and  $\beta^*$  which yield a straight line (called the regression line)  $Y = \alpha^* + \beta^*X$  which minimises a measure of the aggregate distance of the points  $(Y_i, X_i)$ ,  $i = 1, 2, \dots, n$  to that line in X Y space, where Y is measured on the vertical axis. The measure we use is the sum of squared vertical distances which we shall call the Error Sum of Squares (ERSS) so that  $\alpha^*$  and  $\beta^*$  are solutions to the problem:

$$\min_{\alpha^*, \beta^*} ERSS \left( = \sum_{i=1}^n (Y_i - (\alpha^* + \beta^* X_i))^2 \right)$$

The multivariate calculus is employed to solve this problem by setting the partial derivatives of ERSS with respect to  $\alpha^*$  and  $\beta^*$  to zero (called the first order conditions) and solving thus:

$$\frac{\partial ERSS}{\partial \alpha^*} = -2 \sum_{i=1}^n (Y_i - (\alpha^* + \beta^* X_i)) = 0$$
$$\frac{\partial ERSS}{\partial \beta^*} = -2 \sum_{i=1}^n (Y_i - (\alpha^* + \beta^* X_i)) X_i = 0$$

The solutions to which are:

$$\alpha^* = \bar{Y} - \beta^* \bar{X}$$

$$\beta^* = \frac{\sum_{i=1}^n (Y_i - \bar{Y}) X_i}{\sum_{i=1}^n (X_i - \bar{X}) X_i}$$

Note the formula for  $\beta^*$  has many alternative equivalent versions which may be seen by observing that for the numerator:

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

and for the denominator:

$$\sum_{i=1}^n (X_i - \bar{X}) X_i = \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X}) = \sum_{i=1}^n X_i^2 - n \bar{X}^2$$

so that various combinations of numerator and denominator formulae yield 12 alternative representations of  $\beta^*$ .

### **Assumptions in the Ordinary Least Squares model.**

Note that while  $\alpha$ ,  $\beta$  and  $\varepsilon_i$ ,  $i = 1, \dots, n$  are fundamentally unobservable we only concern ourselves with estimating  $\alpha$  and  $\beta$  which define the relationship between Y and X. The  $\varepsilon_i$ ,  $i = 1, \dots, n$  are considered “errors” which accommodate all the other influences on Y not accounted for in  $\alpha + \beta X$  as such we assume them to be random and to obey the following four assumptions:

**1)  $E(\varepsilon_i) = 0$  for all  $i$ .**

This assumption really says that the average or net effect of all the other influences on Y not accounted for in  $\alpha + \beta X$  is constant and zero for each observation ( $i=1, \dots, n$ ).

**2)  $V(\varepsilon_i) = \sigma^2 > 0$  for all  $i$ .**

This assumption really says that the variability of the net effect of all the other influences on Y not accounted for in  $\alpha + \beta X$  is constant and non zero for each observation ( $i=1, \dots, n$ ).

This is sometimes referred to as the homoskedasticity assumption and is not as innocuous as it seems. For example if Y were the consumption behaviour of an individual and X was their disposable income it says that the scale of variability of the unobserved effects would be the same for both rich and poor individuals.

**3)  $E(\epsilon_i \epsilon_j) = 0$  for all  $i \neq j$ .**

**4)  $E(\epsilon_i X_j) = 0$ . For  $i$  and  $j$**

These last two assumptions derive from the general assumption that the net effects of all the other influences on Y not accounted for in  $\alpha + \beta X$  are independent of the X's and of each other. Again these are not innocuous assumptions, for example if the equation relates to the behaviour of individuals and individuals in the sample are related to one another in some fashion they are unlikely to be true. Similarly if the equation relates to behaviour through time today's random effects are unlikely to be independent of yesterday's effects.

Given these assumptions certain properties of the estimators follow.

### **Unbiasedness.**

Under the above assumptions the ordinary least squares estimators  $\alpha^*$  and  $\beta^*$  are unbiased so that  $E(\alpha^*) = \alpha$  and  $E(\beta^*) = \beta$  which may be demonstrated as follows. From the various formulae for  $\beta^*$  we may write:

$$\beta^* = \frac{\sum_{i=1}^n (X_i - \bar{X}) Y_i}{\sum_{i=1}^n (X_i - \bar{X}) X_i} = \frac{\sum_{i=1}^n (X_i - \bar{X}) (\alpha + \beta X_i + \epsilon_i)}{\sum_{i=1}^n (X_i - \bar{X}) X_i}$$

which, after noting that the sum of deviations from mean is equal to 0, may be written as:

$$\beta^* = \beta + \frac{\sum_{i=1}^n (X_i - \bar{X}) \epsilon_i}{\sum_{i=1}^n (X_i - \bar{X}) X_i}$$

which can readily be seen to have an expectation  $\beta$ . For  $\alpha^*$  note that:

$$\alpha^* = \bar{Y} + \beta^* \bar{X} = \frac{1}{n} \sum_{i=1}^n (\alpha + \beta X_i + \epsilon_i) - \beta^* \bar{X} = \alpha + (\beta - \beta^*) \bar{X} + \bar{\epsilon}$$

since  $E(\beta - \beta^*) = 0$  and given assumption 1  $E(\alpha^*) = \alpha$ .

### The variances of the estimators.

The variance of the estimators  $\alpha^*$  and  $\beta^*$  facilitate inference and confidence interval estimation. The can be derived, given the above assumptions as follows. The variance of any random variable  $W$  is given by  $E[(W-E(W))^2]$  so that for  $\beta^*$  from the above we may write:

$$E\left[\left(\beta^* - E(\beta^*)\right)^2\right] = E\left[\left(\beta + \frac{\sum_{i=1}^n (X_i - \bar{X}) \epsilon_i}{\sum_{i=1}^n (X_i - \bar{X})^2} - \beta\right)^2\right] = E\left[\left(\frac{\sum_{i=1}^n (X_i - \bar{X}) \epsilon_i}{\sum_{i=1}^n (X_i - \bar{X})^2}\right)^2\right]$$

Note from our assumptions that:

$$\begin{aligned} E\left((X_i - \bar{X}) \epsilon_i (X_j - \bar{X}) \epsilon_j\right) &= 0 \quad \text{for all } i \neq j \\ &= \sigma^2 (X_i - \bar{X})^2 \quad \text{for all } i = j \end{aligned}$$

which means that the variance may be written as:

$$E\left[\left(\beta^* - E(\beta^*)\right)^2\right] = \frac{\sum_{i=1}^n (X_i - \bar{X})^2 \sigma^2}{\left(\sum_{i=1}^n (X_i - \bar{X})^2\right)^2} = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

By similar arguments the variance of  $\alpha^*$  is given by:

$$E\left[\left(\alpha^* - E(\alpha^*)\right)^2\right] = \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right)$$

### The Residuals.

The things that are not “explained” by the model, the  $\epsilon_i$ 's are in fact implicitly estimated by the model by  $Y_i - \alpha^* - \beta^* X_i$ . These estimates of the “errors” are frequently referred to as the “Residuals”. They too have properties which reflect the assumptions made regarding the true errors.

1. The residuals sum to 0.

$$\begin{aligned}\sum_{i=1}^n e_i &= \sum_{i=1}^n (Y_i - \alpha^* - \beta^* X_i) = \sum_{i=1}^n (Y_i - (\bar{Y} - \beta^* \bar{X}) - \beta^* X_i) \\ &= \sum_{i=1}^n (Y_i - \bar{Y}) - \beta^* \sum_{i=1}^n (X_i - \bar{X}) = 0\end{aligned}$$

This mimics the idea that the errors have an expectation 0.

2. The residuals are orthogonal to the X's.

$$\begin{aligned}\sum_{i=1}^n e_i X_i &= \sum_{i=1}^n (Y_i - \alpha^* - \beta^* X_i) X_i = \sum_{i=1}^n (Y_i - (\bar{Y} - \beta^* \bar{X}) - \beta^* X_i) X_i \\ &= \sum_{i=1}^n (Y_i - \bar{Y}) X_i - \beta^* \sum_{i=1}^n (X_i - \bar{X}) X_i = \sum_{i=1}^n (Y_i - \bar{Y}) X_i - \frac{\sum_{i=1}^n (Y_i - \bar{Y}) X_i}{\sum_{i=1}^n (X_i - \bar{X}) X_i} \sum_{i=1}^n (X_i - \bar{X}) X_i = 0\end{aligned}$$

Which mimics the idea that the errors are independent of the  $X_i$  's.

### **Inference and Confidence Intervals.**

Noting that  $\beta^*$  may be written as:

$$\beta^* = \beta + \frac{\sum_{i=1}^n (X_i - \bar{X}) \epsilon_i}{\sum_{i=1}^n (X_i - \bar{X}) X_i}$$

it can readily be seen that the estimator is a linear function of the errors  $\epsilon_i$ ,  $i = 1, \dots, n$  so that if they were normally distributed then so would  $\beta^*$  be, in fact based upon the foregoing it is distributed as:

$$\beta^* \sim N \left( \beta, \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right)$$

in a similar fashion  $\alpha^*$  will be distributed as:

$$\alpha^* \sim N \left( \alpha, \sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right) \right)$$

so that inference and confidence interval computations can proceed accordingly.

### An Estimator for $\sigma^2$ .

As usual it is extremely rare for the value of the variance to be known so that generally it has to be estimated. The estimator  $\sigma^{2*}$  is given by problem:

$$\sigma^{2*} = \frac{ERSS}{n-2} = \frac{\sum_{i=1}^n (Y_i - (\alpha^* + \beta^* X_i))^2}{n-2}$$

which is most conveniently calculated by noting that:

$$ERSS = \sum_{i=1}^n (Y_i - \bar{Y})^2 - \beta^{*2} \sum_{i=1}^n (X_i - \bar{X})^2$$

employing this estimate means that:

$$\frac{(\beta^* - \beta)}{\sqrt{\frac{\sigma^{*2}}{\sum_{i=1}^n (X_i - \bar{X})^2}}} \sim t(n-2)$$

and similarly:

$$\frac{(\alpha^* - \alpha)}{\sqrt{\sigma^2 \left( \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right)}} \sim t(n-2)$$

so inference and confidence interval calculations can be conducted accordingly.

## **R<sup>2</sup> and all that.**

A very common instrument for measuring the degree of explanatory power in an equation is the R<sup>2</sup> statistic which measures the proportion of the variability in the Y<sub>i</sub>'s that is attributable to the X<sub>i</sub>'s. Given the amount of variability in the Y<sub>i</sub>'s attributable to the errors is given by the error sum of squares (ERSS), the proportion of the variability of the Y<sub>i</sub>'s attributable to the X<sub>i</sub>'s is given by:

$$R^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2 - ERSS}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \frac{\beta^{*2} \sum_{i=1}^n (X_i - \bar{X})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

where the last equality is obtained by substituting in the formula for ERSS. Noting that 1-R<sup>2</sup> is:

$$1 - R^2 = \frac{ERSS}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

it follows that:

$$\frac{1 - R^2}{R^2} \cdot (n - 2) = \frac{\beta^{*2} \sum_{i=1}^n (X_i - \bar{X})^2}{\frac{ERSS}{n - 2}} = \left( \frac{\beta^*}{st\ dev(\beta^*)} \right)^2$$

which is the square of the classic “t” statistic for β which will have an F(1,n-2) distribution.

## **Multiple Regression.**

Here we consider the extension of the simple regression technique to analysing the linear relationship between K+1 variables Y and X<sub>1</sub>, X<sub>2</sub>, ..., X<sub>K</sub>. We have n K+1-tuples of observations (Y<sub>i</sub> X<sub>i1</sub>, X<sub>i2</sub>, ..., X<sub>iK</sub>) i = 1, 2, ..., n on the relationship which, because it is not exact, we shall write as:

$$Y_i = \alpha + \sum_{k=1}^K \beta_k X_{ik} + \epsilon_i \quad i = 1, \dots, n$$

In this relationship  $\alpha$ ,  $\beta_k$ ,  $k = 1, \dots, K$  and  $\varepsilon_i$   $i = 1, \dots, n$  are fundamentally unobservable and we would like to estimate the  $\alpha$  and  $\beta_k$ . Essentially this deals with the case where  $Y$  is described by more than one variable  $X$ . We require one addition to our 4 assumptions, namely  $n > K+1$ .

In the same way that estimators for  $\alpha$ ,  $\beta$  were developed in the simple regression case (by minimising ERSS) estimators for  $\alpha$ ,  $\beta_k$ ,  $k = 1, \dots, K$  (denoted  $\alpha^*$ ,  $\beta_k^*$ ,  $k = 1, \dots, K$ ) can be developed by minimizing with respect to  $\alpha^*$ ,  $\beta_k^*$ ,  $k = 1, \dots, K$  the multivariate version of the error sum of squares given by:

$$ERSS = \sum_{i=1}^n \left( Y_i - \alpha^* + \sum_{k=1}^K \beta_k^* X_{ik} \right)^2$$

The formulae for these estimates and similarly formulae for  $\text{Var}(\alpha^*)$  and  $\text{Var}(\beta_k^*)$ ,  $k = 1, \dots, K$  are complicated and will not be given here however these estimates are readily calculated using available software packages and just as in the simple regression case inference and confidence intervals may be pursued by noting that:

$$\frac{(\beta_k^* - \beta_k)}{\sqrt{\text{var}(\beta_k^*)}} \sim t(n - K - 1)$$

and similarly:

$$\frac{(\alpha_k^* - \alpha_k)}{\sqrt{\text{var}(\alpha_k^*)}} \sim t(n - k - 1)$$

so that inference and confidence interval calculations can be conducted accordingly.

$R^2$  in the Multivariate Case.

In the multivariate case  $R^2$  can be calculated in exactly the same fashion as:

$$R^2 = 1 - \frac{ERSS}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$



however in this case if we wish to use it in the context of a joint test of the significance of all of the  $X_{ik}$ 's we have to use:

$$\frac{R^2}{1-R^2} \cdot \frac{(n-K-1)}{K} \sim F(K, n-K-1)$$

which provides us with a one sided upper tailed test statistic for the significance of the K explanatory variables  $X_{ik}$ .