

## LECTURE 09

### SIMPLE REGRESSION MODEL - I

Outline of today's lecture:

I. Simple Linear Regression Model .....	1
II. Linearity Issue .....	2
III. Stochastic Nature of Linear Regression Model .....	3
IV. Assumptions of Classical Linear Model.....	4
V. Gauss Markov Theorem.....	8
VI. Normality Assumption.....	8

### I. Simple Linear Regression Model

- In econometrics we deal exclusively with *stochastic* relations.
  - Stochastic is a term for *random* or *uncertain*.
  - It is the opposite of *deterministic*.
- The simplest form of stochastic relation between two variables X and Y is called a *simple linear regression* model.
  - $Y_t = \beta_0 + \beta_1 X_t + u_t \quad t=1, \dots, n$

where

- Y is *dependent variable*
- X is *independent* (explanatory) variable<sup>1</sup>
- u is the stochastic *disturbance term*, or *error term*.
- $\beta_0$  and  $\beta_1$  are the regression *parameters* which are *unknown*.
- Subscript t refers to the t<sup>th</sup> observation.

---

<sup>1</sup> The terms *regressor*, *regressand* and *covariate* are also used.

- The values of the variables  $X$  and  $Y$  are observable.
- However, the values of  $u$  are not observable.
- Observations on  $X$  and  $Y$  can be made over time, in which case we speak of having “*time series*” data.
  - For example we may have data on Turkey’s GDP over 30 years (data collected over a period of time).
- Or they can be made over
  - individuals, or groups of individuals,
  - firms, or group of firms,
  - countries, or group of countries,
  - objects,
  - geographical areas, etc.,

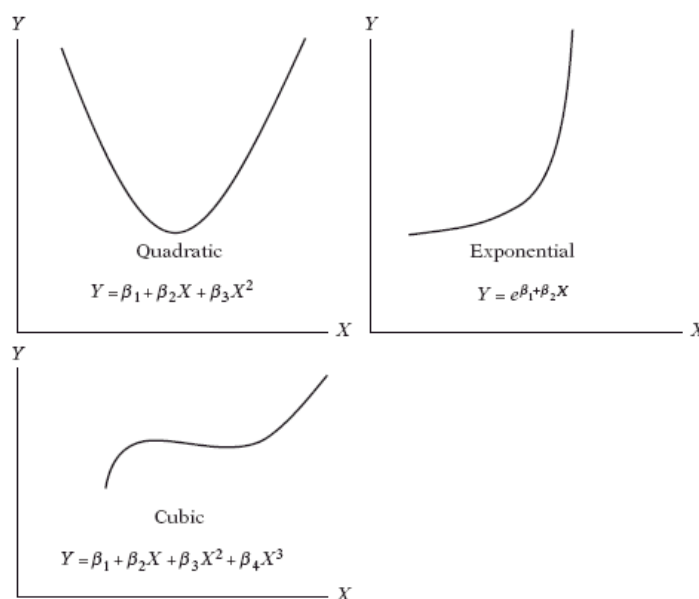
in which we speak of having “*cross section data*”.

- For example we may have data on several countries’ GDP for one year (data collected at one point in time).
- Hence, the subscript  $t$  may refer to the  $t^{\text{th}}$  year (quarter, month, day, etc) or to the  $t^{\text{th}}$  individual or group (such as countries; Turkey, Germany, France, USA, Japan, etc..).
- Data of both kinds can be combined to obtain “*pooled times and cross-section data*”.
  - For example we may have data on 20 countries’ GDP over 30 years.

## II. Linearity Issue

- $Y_t = \beta_0 + \beta_1^3 X_t + u_t$  is nonlinear in the parameter  $\beta_1$ .

- This model is an example of a nonlinear (in the parameter) regression model.



Linear-in-parameter functions.

LINEAR REGRESSION MODELS		
Model linear in parameters?	Model linear in variables?	
	Yes	No
Yes	LRM	LRM
No	NLRM	NLRM

Note: LRM = linear regression model  
 NLRM = nonlinear regression model

### III. Stochastic Nature of Linear Regression Model

- The stochastic nature of the regression model implies that for every value of  $X$  there is a whole probability distribution of values of  $Y$ .
- This means that the value of  $Y$  can *never* be forecast exactly.
- *Uncertainty* concerning  $Y$  arises because of the presence of the stochastic disturbance  $u$  which, being random, imparts randomness to  $Y$ .
  - For example, consider a production function of a firm.

- Suppose that output depends in some specified way on the quantity of labor (L) in accordance with firm's technology.

$$\text{Output}_t = \beta_0 + \beta_1 \text{Labor}_t + u_t$$

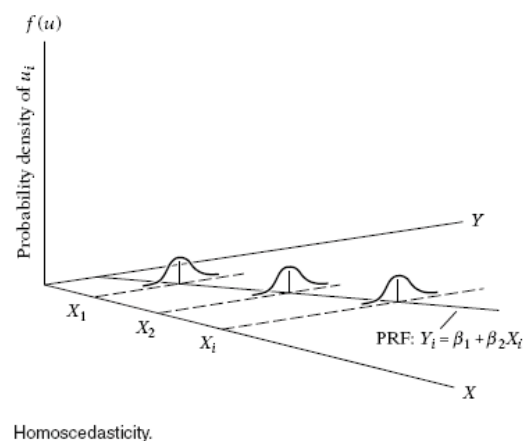
- This function may apply in the short run when the quantities of other factors are fixed.
- In general, the same quantity of labor will lead to different quantities of output because of variations in *weather*, *human performance*, *frequency of machine breakdowns*, and many other factors.
- Output, which is the dependent variable in this case, will depend not only on the quantity of labor input, but also on the large number of random causes,
  - *which we summarize in the form of the stochastic disturbance (u).*

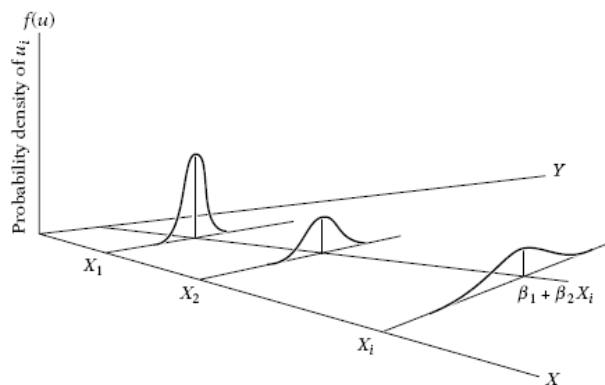
#### IV. Assumptions of Classical Linear Model

- **Assumption 1:** Linear regression model [the regression model is linear in the parameters]
- **Assumption 2:**  $X$  values are fixed in repeated sampling

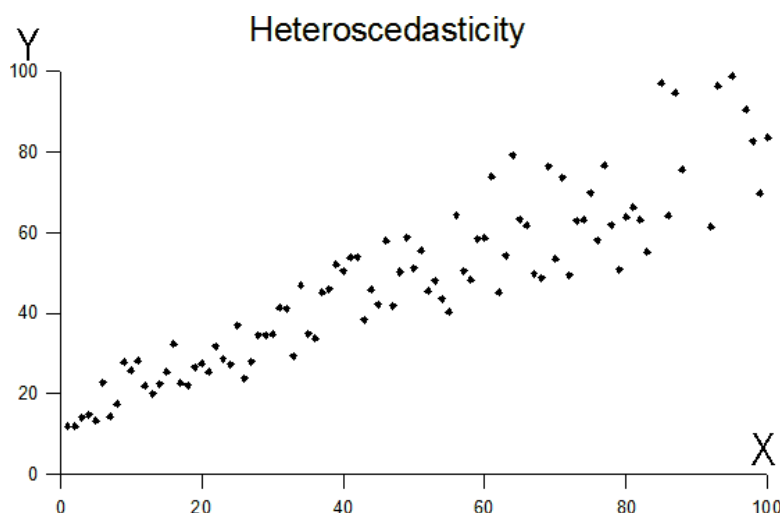
This means that in taking a large number of samples on  $Y$  and  $X$ , the  $X_t$  values are the same in all samples, but the  $u_t$  values differ from sample to sample, and so of course do the values of  $Y_t$ . For example assume that every day in market we choose the same prices  $X_1, X_2, \dots, X_t$ , and we record the quantities of  $Y_t$ 's sold each day at these prices. The  $X$ 's do not vary; they are a set of fixed values, while the  $Y_t$ 's vary for each day due to different random influences.

- **Assumption 3:** Zero mean value of disturbance  $u_t$  [ $E(u_t)=0$ ]
  - This condition is that the expected value of the disturbance term in any observation should be 0.
  - Sometimes it will be positive, sometimes negative, but it should not have a systematic tendency in either direction.
  - This assumption says that the factors not explicitly included in the model, and therefore subsumed in  $u_t$ , do not systematically affect the mean value of Y.
  
- **Assumption 4:** *Homoscedasticity* [ $Var(u_t)=\sigma^2$ ].
  - Population variance of the disturbance term ( $u_t$ ) should be constant for all observations (Homoscedasticity).
  - If this condition is not satisfied, the OLS regression coefficients will be *inefficient*.
  - *In terms of our production example, this assumption implies that the variation in output is the same whether the quantity of labor is 20, 100, or any other number of units.*





Heteroscedasticity.



- **Assumption 5:** No autocorrelation [correlation between any  $u_t$  and  $u_s$  ( $t \neq s$ ) is zero]
  - This condition states that there should be no systematic association between the values of the disturbance term in any two observations.
  - If this condition is not satisfied, OLS will again give inefficient estimates.
  - *In terms of our production example, this assumption implies that output is higher than expected today should not lead to higher (or lower) than expected output tomorrow.*
  - Recall that *correlation* between X and Y is given by:

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

where  $\sigma_X$  and  $\sigma_Y$  are standard deviations of X and Y, respectively. Therefore, if the correlation between X and Y is zero ( $\rho = 0$ ), it implies that  $Cov(X, Y) = 0$ . As a result, the autocorrelation implies that  $Cov(u_t, u_s) = 0$  where  $t \neq s$ .<sup>2</sup>

- Note that covariance between X and Y is given by

$$Cov(X, Y) = E \{ [X - E(X)][Y - E(Y)] \}$$

Therefore, no autocorrelation implies that

$$Cov(u_t, u_s) = E \{ [u_t - E(u_t)][u_s - E(u_s)] \} = 0$$

From *Assumption 3*, we know that  $E(u_t) = 0$  and  $E(u_s) = 0$ . Thus, autocorrelation assumption implies that

$$Cov(u_t, u_s) = E \left\{ \left[ \begin{array}{c} u_t - \underbrace{E(u_t)}_0 \\ \left[ \begin{array}{c} u_s - \underbrace{E(u_s)}_0 \end{array} \right] \end{array} \right] \right\} = 0 \Rightarrow$$

$$E(u_t u_s) = 0$$

- **Assumption 6:** Zero covariance between  $u_t$  and  $X_t$  [ $E(u_t X_t) = 0$ ]
- **Assumption 7:** The number of observations  $n$  must be greater than the number of parameters  $k$ .
- **Assumption 8:** Variability in  $X$  values [ $Var(X) > 0$ ]

---

<sup>2</sup> Recall that a zero value of the covariance indicates no linear dependence between X and Y.

- **Assumption 9:** The regression model is correctly specified [no specification bias]
- **Assumption 10:** There is no perfect multicollinearity [there are no perfect linear relationships among explanatory variables]

## V. Gauss Markov Theorem

**Definition** Given the assumptions of the CLRM<sup>3</sup>, the OLS estimators, in the class of unbiased linear estimators, have minimum variance<sup>4</sup>, that is they are *Best*<sup>5</sup> *Linear Unbiased Estimators (BLUE)*.

## VI. Normality Assumption

- In addition to the Gauss–Markov conditions, one usually assumes that the disturbance term is normally distributed.
- Reason is that if  $u$  is normally distributed, so will be the regression coefficients, and this will be useful to us later when we come to the business of performing tests of hypotheses and constructing confidence intervals for  $\beta_1$  and  $\beta_2$  using the regression results.
- Justification for the assumption depends on the Central Limit Theorem.

---

<sup>3</sup> See Handout 10 for the details of CLRM.

<sup>4</sup> In other words, they are efficient estimators.

<sup>5</sup> In other words, they are efficient or they have minimum variance.