

Contents

1	Introduction	1
1.1	Definition of Econometrics	1
1.2	Economic Model Building	1
1.3	Scope and Objective of Econometrics	3
1.4	Steps in Econometric Analysis	4
1.5	Example on Simultaneous Equations Model	5
2	Review of Regression and Problems in Econometrics	9
2.1	Simple linear regression	9
2.2	Multiple Linear Regressions	41
2.3	Hetroscedasticity	59
2.4	Autocorrelation	80
2.5	Multicollinearity	115
2.6	Errors in Variables	124
2.7	2.7Lagged variables	127
2.8	2.8 Further Model Peculiarities	127
3	Non-linear Models	154
3.1	Introduction	154
3.2	Intrinsically Linear Models	154
3.3	Intrinsically Non-linear Models	155
3.4	Estimation of Intrinsically non-linear Models	162
4	Simultaneous Equations Models (SEM)	166
4.1	Introduction to SEM	166
4.2	The identification problem	173
4.3	Conditions for identification	174
4.4	Estimation of simultaneous equations models: ILS, 2SLS	178



ECONOMETRICS USING R

By

TEDROS GEBREGERGS

MEKELLE UNIVERSITY

STATISTICS DEPARTMENT

May 2021

Course Guide Book

Course Titles/Code: Econometrics (Stat3061)

Credit: 5 EtCTS

Module title/code: Statistical Modeling II (Stat-M3061)

Course Type: Core

Preface

This module is prepared to serve as a teaching material for the course Econometrics mainly given to undergraduate statistics students. In addition, it can also be used by different instructors and students who work with this course generally, and undergraduate economics program particularly.

This module is organized in four chapters. Chapter one gives general introduction about the course econometrics. The next three chapters deal with theoretical and applied econometric models. Hence, it describes the nature of basic econometrics model, they also provide a practical related applications or solutions.

At the end of each chapter, the writer put self-test exercises. This exercise covers a wide variety of questions that help the learners to relate the theoretical concepts to practical applications.

Besides the theoretical and practical matters covered in this module, the course needs dealing with a practical dataset. To this end, students should have access to a platform where they get or manipulate such dataset. The writer has developed an independent dataset R package extension called “ecostatmomona”. Students and instructors can make an access to the package as follows:

```
devtools::install_github("hahustat/ ecostatmomona")  
library (ecostatmomona)
```

Furthermore, they can also have a supportive interactive online application program at www.github/hahustat/econometrics . Using the following detailed address and procedure will help its users to be at ease.

```
devtools::install_github("hahustat/econometricsR")  
library(swirl)
```

Contents

1 Introduction

1.1 Definition of Econometrics

Literally interpreted, econometrics means “economic measurement”. Econometrics deals with the measurement of economic relationships. It is an integration of economics, mathematical economics and statistics with an objective to provide empirical values to the parameters of economic relationships.

The relationships of economic theories are usually expressed in mathematical forms and combined with empirical economics. For example, to express a relationship between income and consumption, we may write

$$CONSUMPTION = f(INCOME)$$

which says that the level of consumption is some function, $f(*)$, of income.

The econometrics methods are used to obtain the values of parameters which are essentially the coefficients of the mathematical form of the economic relationships. The statistical methods which help in explaining the economic phenomenon are adapted as econometric methods. The econometric relationships depict the random behavior of economic relationships which are generally not considered in economics and mathematical formulations.

Econometrics aims at giving empirical content to economic relationships. The exciting thing about econometrics is its concern for verifying or refuting an economic theory. These economic theory or hypotheses are testable with economic data.

1.2 Economic Model Building

The first task an econometrician faces is that of **formulating an econometric model** but it is better to define **What is model?**

- A model is a simplified representation of a real-world process.
- The goal of a model is to provide a simple low-dimensional summary of a dataset.

- The goal of a model is not to uncover truth, but to discover a simple approximation that is still useful.

Model should be representative in the sense that it should contain the salient features of the phenomena under study. In general, one of the objectives in modeling is to have a simple model to explain a complex phenomenon. For instance, ‘the demand for oranges depends on the price of oranges’ is a simplified representation since there are a host of other variables that one can think of that determine the demand for oranges. These include:

- Income of consumers
- An increase in diet consciousness (e.g. drinking coffee causes cancer; so better switch to orange juice)
- Increase or decrease in the price of substitutes (e.g. that of apple)

However, there is no end to this stream of other variables! Many have argued in favour of simplicity since simple models are easier:

- To understand
- To communicate
- To test empirically with data

The choice of a simple model to explain complex real-world phenomena leads to two criticisms:

- The model is oversimplified
- The assumptions are unrealistic

For instance, to say that the demand for oranges depends on only the price of oranges is both an oversimplification and an unrealistic assumption.

- To the criticism of oversimplification, many have argued that it is better to start with a simplified model and progressively construct more complicated models.
- As to the criticism of unrealistic assumptions, the relevant question is whether they are sufficiently good approximations for the purpose at hand or not

In practice, generally, all the variables which the experimenter thinks are relevant to explain the phenomenon are included in the model. Rest of the variables are dumped in a basket called “disturbances” where the disturbances are random variables. This is the main difference between

economic modeling and **econometric modeling**. This is also the main difference between mathematical modeling and statistical modeling. The mathematical modeling is exact in nature, whereas the statistical modeling contains a stochastic term also.

An economic model is a set of assumptions that describes the behavior of an economy, or more generally, a phenomenon.

An econometric model consists of

- ✓ a set of equations describing the behaviour. These equations are derived from the economic model and have two parts – observed variables and disturbances.
- ✓ a statement about the errors in the observed values of variables.
- ✓ a specification of the probability distribution of disturbances.

1.3 Scope and Objective of Econometrics

The three main objective econometrics are as follows:

1. Formulation and specification of econometric models: The economic models are formulated in an empirically testable form. Several econometric models can be derived from an economic model. Such models differ due to different choice of functional form, specification of the stochastic structure of the variables etc.

2. Estimation and testing of models: The models are estimated on the basis of the observed set of data and are tested for their suitability. This is the part of the statistical inference of the modelling. Various estimation procedures are used to know the numerical values of the unknown parameters of the model. Based on various formulations of statistical models, a suitable and appropriate model is selected.

3. Use of models: The obtained models are used for forecasting and policy formulation, which is an essential part in any policy decision. Such forecasts help the policymakers to judge the goodness of the fitted model and take necessary measures in order to re-adjust the relevant economic variables. The scope of econometrics ranges its theoretical and application according to types of data. The various types of data is used in the estimation of the econometric model includes time series data, cross-section data, panel data, spatial data and others.

1.4 Steps in Econometric Analysis

How do econometricians proceed in their analysis of an economic problem? That is, what is their methodology?

An econometric methodology proceeds along the following lines:

- Formulation of theory or hypothesis
- Specification of economic (mathematical) model,
- Specification of econometric model,
- Collecting data,
- Estimation of parameters,
- Hypothesis tests,
- Forecasting/Prediction),
- Evaluation of results for policy analysis or decision making

Description of the steps

1. We begin with an **economic model** which is a set of **assumptions** that describes the behaviour of an **economic phenomenon**.
2. Formulation (specification) of an economic model: a **set of equations** derived from the economic model that involve some **observed variables** and some **‘disturbances’**.
3. Collection of **relevant data** on variables implied by the econometric model.
4. **Estimation of model parameters** using mathematical statistics and probability theory.
5. We **conduct tests** to verify whether:
 - The specification of the model is correct
 - Model assumptions are valid
6. Based on step:
 - If the **model failed to pass the specification testing and diagnostic checking step**, then **revise** .
 - If the model passes , then one has to **proceed** with testing any hypothesis of interest (e.g. which of the explanatory variables significantly affect the response variable?).
7. We use the estimated model for **predictions and policy**

Example about economic model Vs econometric model

1) consumption (**c**) a linear function of income (**i**),

$$c = f(i) = \beta_0 + \beta_1 income \implies \text{Economic model (Mathematical Model)}$$

$$c = f(i) = \beta_0 + \beta_1 income + \epsilon \implies \text{Econometric mode}$$

2) Job Training and Worker Productivity

$$wage = f(educ, exper, training)$$

where

wage : hourly wage (in dollars)

educ : level of education (in years)

exper : level of workforce experience (in years)

training : weeks spent in job training.

$$wage = \beta_0 + \beta_1 educ + \beta_2 exper + \beta_3 training : \implies \text{Economic Model}$$

$$wage = \beta_0 + \beta_1 educ + \beta_2 exper + \beta_3 training + \epsilon : \implies \text{Econometric Model}$$

1.5 Example on Simultaneous Equations Model

Economists formulate models for consumption, production, investment, money demand and money supply, labor demand and labor supply to attempt to explain the workings of the economy. These behavioral equations are estimated equation by equation or jointly as a system of equations. These are known as simultaneous equations models.

Example 1: Keynesian model of income determination

Consider the simple Keynesian model of income determination:

$$\text{Consumption function : } C_t = \alpha + \beta Y_t + u_t \quad 0 < \beta < 1 \quad t = 1, 2, \dots, T \quad (1.1)$$

$$Y_t = C_t + I_t \quad (1.2)$$

where

C=consumption expenditure

Y=income

I =investment (assumed exogenous)

S=savings

t =time

u=stochastic disturbance term

α and β =parameters

This is a system of two simultaneous equations, also known as structural equations with the second equation being an identity.

The parameter β is known as the marginal propensity to consume(MPC) (the amount of extra consumption expenditure resulting from an extra dollar of income). From economic theory, β is expected to lie between 0 and 1. Equation (1.1) is the (stochastic) consumption function; and (1.2) is the national income identity, signifying that total income is equal to total consumption expenditure plus total investment expenditure, it being understood that total investment expenditure is equal to total savings. Diagrammatically, we have Figure 1.1.

From the postulated consumption function and Figure 1.1 it is clear that C and Y are interdependent and that Y_t in (1.1) is not expected to be independent of the disturbance term because when u_t shifts (because of a variety of factors subsumed in the error term), then the consumption function also shifts, which, in turn, affects Y_t . Therefore, once again the classical least squares method is inapplicable to (1.1). If applied, the estimators thus obtained will be inconsistent, as we shall show in chapter 4.

Example 2: Demand and supply model

As is well known, the price P of a commodity and the quantity Q sold are determined by the

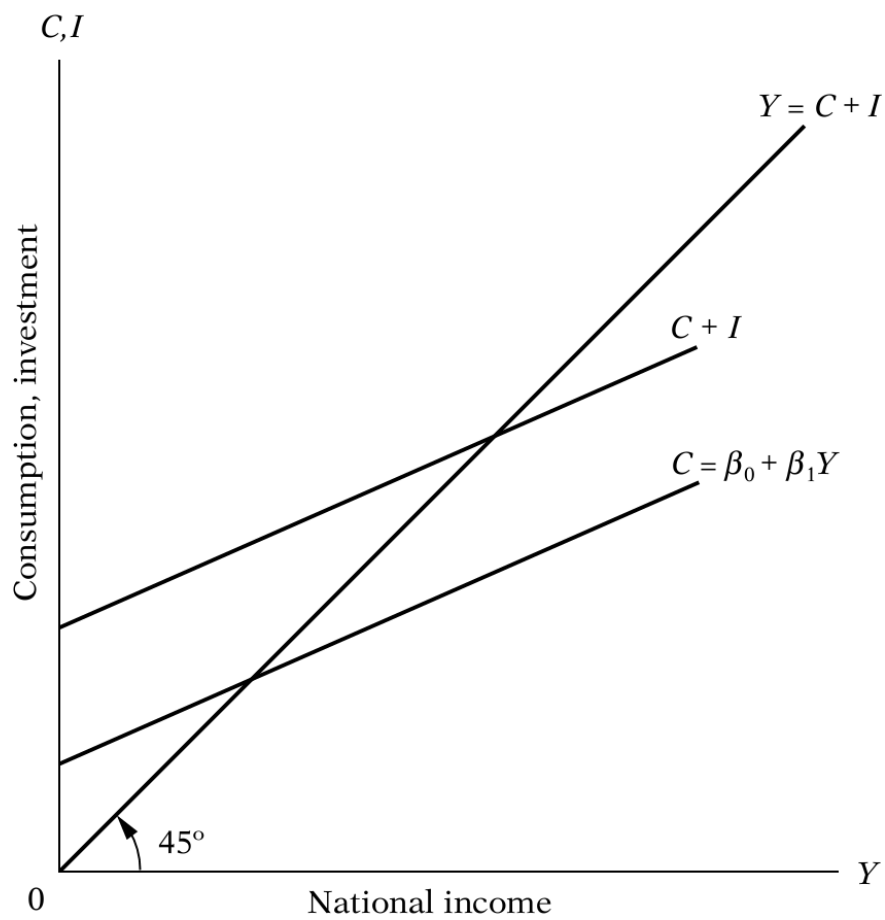


Figure 1.1: Keynesian model of income determination.

intersection of the demand and supply curves for that commodity. Thus, assuming for simplicity that the demand and supply curves are linear and adding the stochastic disturbance terms u_{1t} and u_{2t} , we may write the empirical demand and supply functions as

$$\text{Demand function : } Q_t^d = \alpha_0 + \alpha_1 P_t + u_{1t} \quad \alpha_1 < 0 \quad (1.3)$$

$$\text{Supply function : } Q_t^s = \beta_0 + \beta_1 P_t + u_{2t} \quad \beta_1 > 0 \quad (1.4)$$

$$\text{Equilibrium condition : } Q_t^d = Q_t^s$$

where $Q^d = \text{quantity demanded}$

$Q^s = \text{quantity supplied}$

$t = \text{time}$

and the α 's and β 's are the parameters. A priori, α_1 is expected to be negative (downward-sloping demand curve), and β_1 is expected to be positive (upward-sloping supply curve). Now it is not too difficult to see that P and Q are jointly dependent variables. If, for example, u_{1t} in ((1.3) changes because of changes in other variables affecting Q_t^d (such as in-come, wealth, and tastes), the demand curve will shift upward if u_{1t} is positive and downward if u_{1t} is negative. These shifts are shown in Figure (1.2. As the figure shows, a shift in the demand curve changes both P and Q . Similarly, a change in u_{2t} (because of strikes, weather, import or export restrictions, etc.) will shift the supply curve, again affecting both P and Q . Because of this simultaneous dependence between Q and P , u_{1t} and P_t in ((1.3) and u_{2t} and P_t in ((1.4) cannot be independent. Therefore, a regression of Q on P as in ((1.3) would violate an important assumption of the classical linear regression model, namely, the assumption of no correlation between the explanatory variable(s) and the disturbance term.

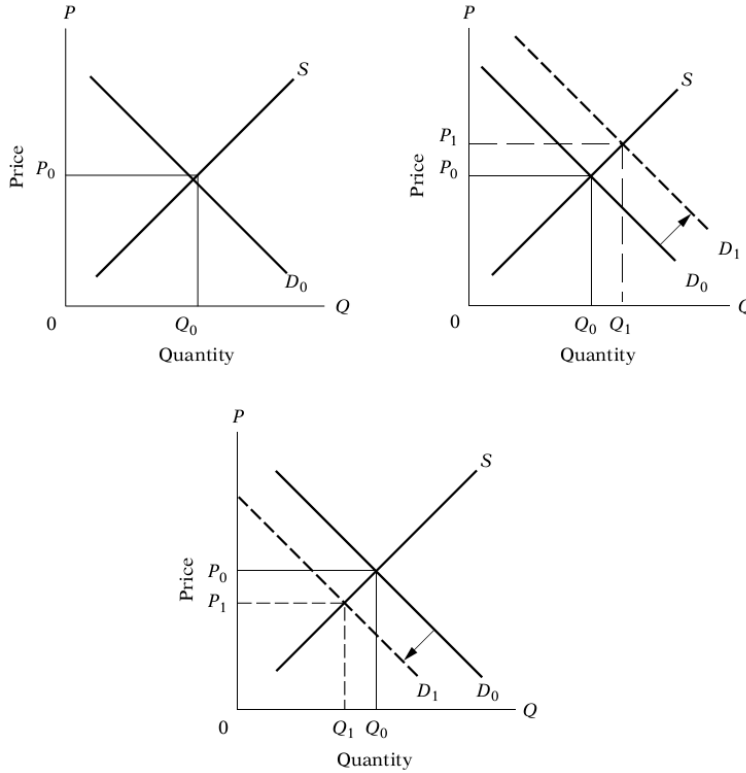


Figure 1.2: Interdependence of price and quantity.

2 Review of Regression and Problems in Econometrics

2.1 Simple linear regression

Introduction

Regression analysis is one of the most commonly used tools in econometric work. Regression analysis is concerned with describing and evaluating the relationship between a given variable (often called the **dependent variable**) and one or more variables which are assumed to influence the given variable (often called **independent** or **explanatory variables**).

Regression analysis is the method to discover the relationship between one or more response variables (also called dependent variables, explained variables, predicted variables, or regressands, usually denoted by y) and the predictors (also called independent variables, explanatory variables, control variables, or regressors, usually denoted by x_1, x_2, \dots, x_k).

Table 2.1: Weekly Family Income X

X	X80	X100	X120	X140	X160	X180	X200	X220	X240	X260
2	55	65	79	80	102	110	120	135	137	150
3	60	70	84	93	107	115	136	137	145	152
4	65	74	90	95	110	120	140	140	155	175
5	70	80	94	103	116	130	144	152	165	178
6	75	85	98	108	118	135	145	157	175	180
7	NA	88	NA	113	125	140	NA	160	189	185
8	NA	NA	NA	115	NA	NA	NA	162	NA	191
Total	25	462	445	707	678	750	685	1043	966	1211
E(Y X)	65	77	89	101	113	125	137	149	161	173

2.1.1 Population Regression Function

```
demo1<-read.csv("PRF.csv",header = T,sep = ",")
demo2<-read.csv("demo2.csv",header = T,sep = ",")
```

The data in Table 2.1 refers to a total population of 60 families in a hypothetical community and their weekly income (X)(columns) and weekly consumption expenditure (Y)(first 7 rows), both in dollars. The 60 families are divided into 10 income groups (from \$80 to \$260) and the weekly expenditures of each family in the various groups are as shown in the table (or read `demo2.csv` please). Therefore, we have 10 fixed values of X and the corresponding Y values against each of the X values; so to speak, there are 10 Y sub populations and look Figure 2.3.

There is considerable variation in weekly consumption expenditure in each income group, which can be seen clearly from Figure 2.3. Despite the variability of weekly consumption expenditure within each income group, on the average, weekly consumption expenditure increases as income increases. corresponding to the weekly income level of \$80, the mean consumption expenditure is \$65, while corresponding to the income level of \$200, it is \$137.

In all we have 10 mean values for the 10 sub populations of Y. We call these mean values **conditional expected values**, as they depend on the **given values** of the (conditioning) variable X (see also Figure 2.3). Symbolically, we denote them as $E(Y|X)$, which is read as the expected value of Y given the value of X.

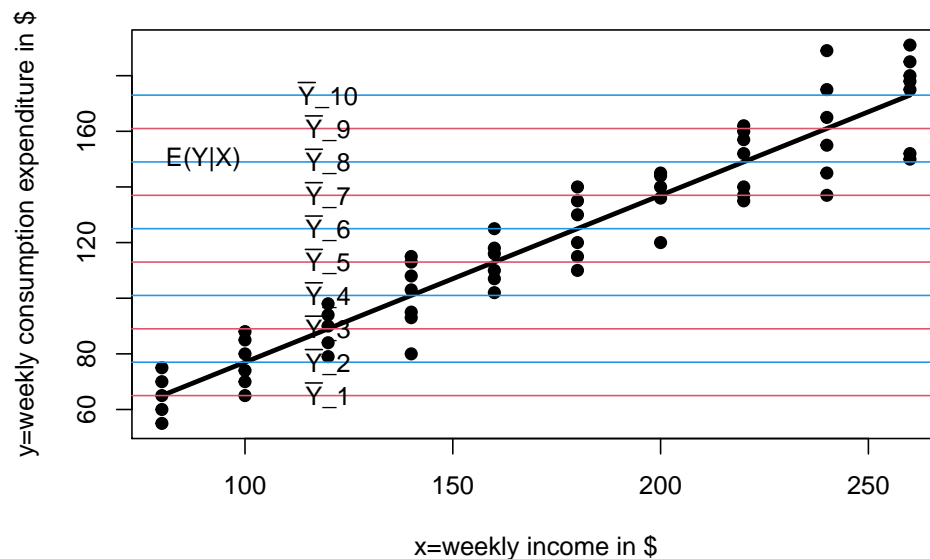


Figure 2.3: Conditional distribution of expenditure for various levels of income

Unconditional expected value of weekly consumption expenditure, $E(Y)$, is obtained adding the weekly consumption expenditures for all the 60 families in the population and dividing this number by 60 ($7272/60 = 121.2$), we get the number \$121.20.

When we ask the question, “What is the expected value of weekly consumption expenditure of a family,” we get the answer \$121.20 (the unconditional mean). But if we ask the question, “What is the expected value of weekly consumption expenditure of a family whose monthly income is, say, \$140,” we get the answer \$101 (the conditional mean).

```
mean(demo2$y[demo2$x==140])
```

```
## [1] 101
```

A population regression curve(**population regression line (PRL)**) is simply the locus of the conditional means of the dependent variable for the fixed values of the explanatory variable(s). More simply, it is the curve connecting the means of the sub populations of Y corresponding to the given values of the regressor X. It can be depicted as in Figure 2.4.

Figure 2.4 shows that for each X (i.e., income level) there is a population of Y values (weekly

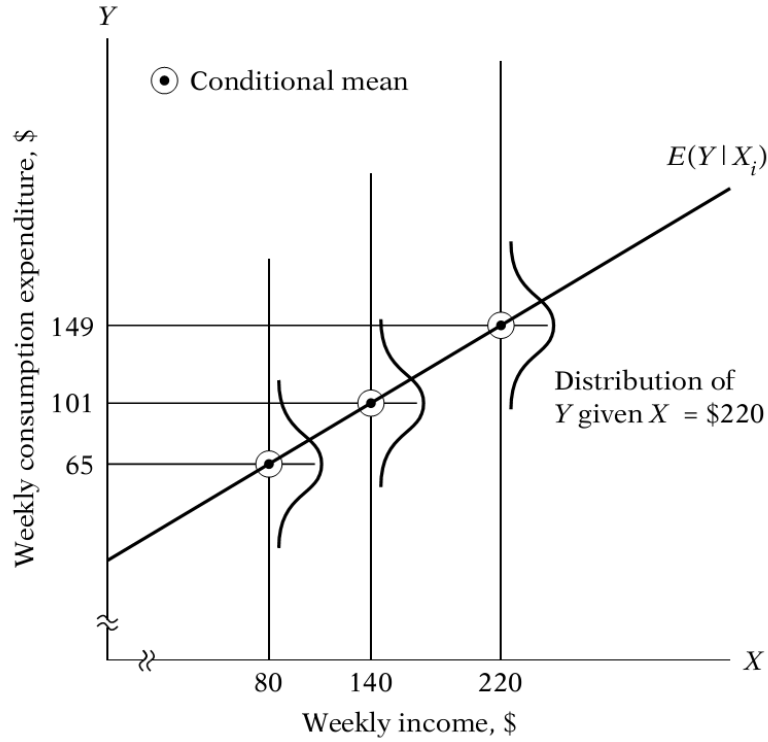


Figure 2.4: Population regression line

consumption expenditures) that are spread around the (conditional) mean of those Y values. For simplicity, we are assuming that these Y values are distributed symmetrically around their respective (conditional) mean values. And the regression line (or curve) passes through these (conditional) mean values.

2.1.1.1 The concept of population regression function(PRF) From the preceding discussions, it is clear that each conditional mean $E(Y|X_i)$ is a function of X_i , where X_i is a given value of X. Symbolically,

$$E(Y|X_i) = f(X_i) \quad (2.5)$$

where $f(X_i)$ denotes some function of the explanatory variable X. Equation (2.5) is known as the conditional expectation function (CEF) or **population regression function (PRF)** or population regression (PR) for short. It states merely that the expected value of the distribution of Y given X_i is functionally related to X_i . In simple terms, it tells how the mean or average response of Y varies

with X.

What form does the function $f(X_i)$ assume? This is an important question because in real situations we do not have the entire population available for examination. The functional form of the PRF is therefore an **empirical question**.

we may assume that the **PRF $E(Y|X_i)$** is a linear function of X_i , say, of type

$$E(Y|X_i) = \alpha + \beta X_i \quad (2.6)$$

where α and β are **unknown but fixed** parameters known as the regression coefficients.

In regression analysis our interest is in estimating the **PRFs** like (2.6), that is, estimating the values of the unknowns α and β on the basis of observations on Y and X.

2.1.1.2 Stochastic specification of PRF It is clear from Figure 2.3, as family income increases, family consumption expenditure on the average increases, too.

Given the income level of X_i , an individual family's consumption expenditure is clustered around the average consumption of all families at that X_i , that is, around its conditional expectation.

Therefore, we can express the deviation of an individual Y_i around its expected value as follows:

$$\begin{aligned} \mathcal{E}_i &= Y_i - E(Y|X_i) \\ \implies Y_i &= E(Y|X_i) + \mathcal{E}_i \end{aligned} \quad (2.7)$$

where the deviation \mathcal{E}_i is an **unobservable random variable** taking positive or negative values.

Technically, \mathcal{E}_i is known as the stochastic disturbance or stochastic error term.

Given the income level, We can say that the expenditure of an individual family can be expressed as the sum of two components:

- 1) $E(Y|X_i)$, which is simply the mean consumption expenditure of all the families with the same level of income. This component is known as the **systematic, or**

deterministic component, and

- 2) \mathcal{E}_i , which is the random, or non systematic component.

If $E(Y|X_i)$ is assumed to be linear in X_i , as in (2.6), equation (2.7) may be written as

$$\begin{aligned} Y_i &= E(Y|X_i) + \mathcal{E}_i \\ &= \alpha + \beta X_i + \mathcal{E}_i \end{aligned} \tag{2.8}$$

Equation (2.8) shows that the consumption expenditure of a family is linearly related to its income plus the disturbance term.

Thus, the individual consumption expenditures, given $X = \$80$ (see Table 2.1), can be expressed as

$$\begin{aligned} Y_1 &= 55 = \alpha + \beta(80) + \mathcal{E}_1 \\ Y_2 &= 60 = \alpha + \beta(80) + \mathcal{E}_2 \\ Y_3 &= 65 = \alpha + \beta(80) + \mathcal{E}_3 \\ Y_4 &= 70 = \alpha + \beta(80) + \mathcal{E}_4 \\ Y_5 &= 75 = \alpha + \beta(80) + \mathcal{E}_5 \end{aligned}$$

Now if we take the expected value of (2.7) on both sides, we obtain

$$\begin{aligned} E(Y_i|X_i) &= E(E(Y|X_i)) + E(\mathcal{E}_i|X_i) \\ &= E(Y|X_i) + E(\mathcal{E}_i|X_i) \\ \implies E(Y_i|X_i) - E(Y|X_i) &= E(\mathcal{E}_i|X_i) \\ \implies E(\mathcal{E}_i|X_i) &= 0 \end{aligned} \tag{2.9}$$

Thus, the assumption that the regression line passes through the conditional means of Y implies that the conditional mean values of \mathcal{E}_i (conditional upon the given X's) are zero.

Example: Suppose the relationship between consumption (Y) and income (X) of households is expressed as:

$$Y_i = 0.6X_i + 120$$

where Y_i =consumption of i^{th} household and X_i =income of i^{th} household

Here, on the basis of income, we can predict consumption. For instance, if the income of a certain household is 1500 Birr, then the estimated consumption will be:

$$\text{consumption} = 0.6(1500) + 120 = 1020 \text{ Birr}$$

Note that since consumption is estimated on the basis of income, consumption is the dependent variable and income is the independent variable.

The error term

Consider the above model: $Y = 0.6X + 120$. This functional relationship is deterministic or exact, that is, given income we can determine the exact consumption of a household. But in reality this rarely happens: different households with the same income are not expected to consume equal amount and this may be due to difference in real wealth or varying tastes, or unforeseen events that induce households to consume more or less.

Thus, we should express the regression model as equation (2.8).

Generally the reasons for including the error term include:

1. **Omitted variables:** a model is a simplification of reality. It is not always possible to include all relevant variables in a functional form . For instance,

- the omission of relevant factors that could influence consumption, other than income, like real wealth or varying tastes, or unforeseen events that induce households to consume more or less.
- we may construct a model relating demand and price of a commodity. But demand is influenced not only by own price: income of consumers, price of substitutes and several other variables also influence it.

The omission of these variables from the model introduces an error.

2. **Measurement error:** inaccuracy in collection and measurement of sample data.
 - households may not report their consumption or income accurately
3. **Sampling error:** Consider a model relating consumption (Y) with income (X) of households. The sample we randomly choose to examine the relationship may turn out to be predominantly poor households. In such cases, our estimation of α and β from this sample may not be as good as that from a balanced sample group.
4. **wrong choice of a linear relationship** between consumption and income, when the **true** relationship may be **nonlinear**.

Note that the size of the error \mathcal{E}_i is not fixed: it is **non-deterministic** or **stochastic** or **probabilistic** in nature. This in turn implies that Y_i is also probabilistic in nature. Thus, the probability distribution of Y_i and its characteristics are determined by the values of X_i and by the probability distribution of \mathcal{E}_i .

Thus, a full specification of a regression model should include a specification of the probability distribution of the disturbance (error) term. This information is given by what we call basic assumptions or assumptions of the classical linear regression model (CLRM).

Consider the model:

$$Y_i = \alpha + \beta X_i + \mathcal{E}_i \quad i = 1, 2, \dots, n.$$

Here the subscript i refers to the i^{th} observation. In the CLRM, Y_i and X_i are observable while \mathcal{E}_i is not. If i refers to some point or period of time, then we speak of time series data. On the other hand, if i refers to the i^{th} individual, object, geographical region, etc., then we speak of cross-sectional data.

2.1.1.3 Assumptions of the classical linear regression model

$$Y_i = \alpha + \beta X_i + \mathcal{E}_i \quad i = 1, 2, \dots, n.$$

Y and X is observable, but not \mathcal{E}

The assumptions (A) are

- **A1:** The true model is: $Y_i = \alpha + \beta X_i + \mathcal{E}_i$
 - the relationship between Y_i and X_i is **linear**, which is linear in the parameters and
 - the deterministic component ($\alpha + \beta X_i$) and the stochastic component (\mathcal{E}_i) are **additive**.
- **A2:** The error terms have **zero mean**: $E(\mathcal{E}_i) = 0$.
 - this assumption tells us that the mean of the Y_i is:

$$E(Y_i) = \alpha + \beta X_i$$

- , which is **non-stochastic**.
 - is needed to insure that **on the average** we are on the **true line**.
- **A3: Homoscedasticity** (error terms have constant variance): $Var(\mathcal{E}_i) = E(\mathcal{E}_i^2) = \sigma^2$ for all i
 - This assumption tells us that **every disturbance** has the **same variance** σ^2 whose value is **unknown**, that is, regardless of whether the X_i are **large or small**, the dispersion of the disturbances is the **same**. For example, the variation in consumption level of low income households is the same as that of high income households.
 - This insures that every observation is **equally reliable**.
- **A4: No error autocorrelation** (the error terms \mathcal{E}_i are statistically independent of each other): $cov(\mathcal{E}_i, \mathcal{E}_j) = E(\mathcal{E}_i \mathcal{E}_j) = 0$ for $i \neq j$.
 - It states that the disturbances are uncorrelated.
 - Knowing the i^{th} disturbance does not tell us anything about the j^{th} disturbance.
 - For example,

- 1) In the consumption case, the **unforeseen** disturbance which caused the i^{th} household to **consume more**, (like a **visit of a relative**), has **nothing** to do with the unforeseen disturbances of any other household. However, this achieves in random sample of households

2) the fact that output is higher than expected today should not lead to a higher (or lower) than expected output tomorrow.

- **A5:** X_i are **deterministic**(non-stochastic): X_i and \mathcal{E}_i are independent for all i, j
- It states that X_i are **not random variables**, and **hence** the probability distribution of \mathcal{E}_i is in no way affected by the X_i .
- $\implies \sum \mathcal{E}_i x_i = 0$
- Fixed in repeated sampling
- **A6: Normality:** \mathcal{E}_i are normally distributed with mean zero and variance σ^2 for all i (often written as: $\mathcal{E}_i \sim N(0, \sigma^2)$).

We need this assumption for parameter estimation purposes and also to make inferences on the basis of the normal (t and F) distribution.

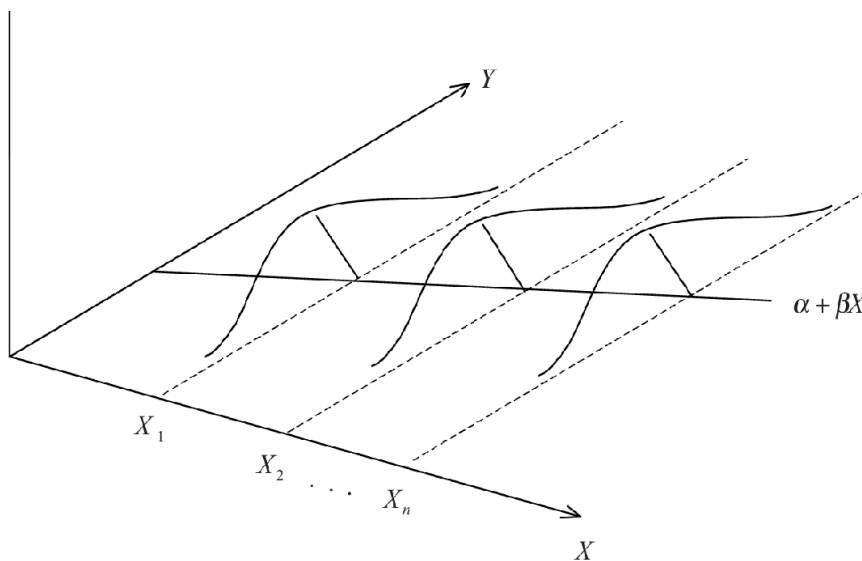


Figure 2.5: Random Disturbances Around the Regression

2.1.2 Estimation methods

2.1.2.1 I. The ordinary least squares (OLS) method of Estimation

- Least Square method is one of the parameter estimation method which minimizes the **sum square of residuals (ESS)**

– Sample Regression Line given by:

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i \quad , i = 1, 2, \dots, n.$$

where α and β are estimated by $\hat{\alpha}$ and $\hat{\beta}$, respectively, and \hat{Y} is the estimated value of Y.

- Residual is the deviation b/n observed and estimated y.

$$\hat{\mathcal{E}}_i = Y_i - \hat{Y}_i \quad , i = 1, 2, \dots, n.$$

- a **good fit** minimizes the error between the **estimated points on the line** and the **actual observed points** , and hence we search which line is?(see Figure 2.6)

Residuals

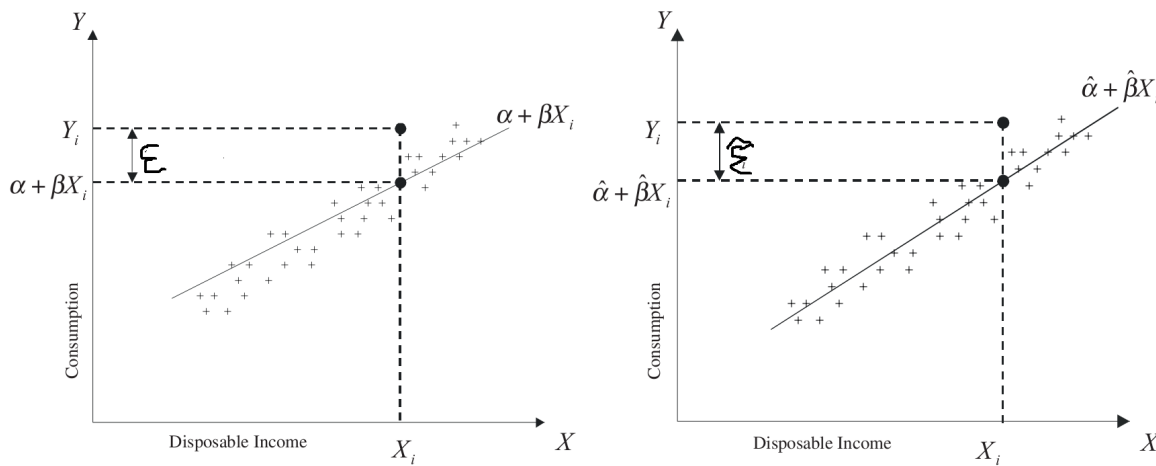


Figure 2.6(a) 'True' Consumption Function

Figure 2.6(b) Estimated Consumption Function

Figure 2.6: True Regression Line and Estimated Line

The sum of squares of the errors (SSE)

$$SSE = \sum \hat{\mathcal{E}}_i^2 = \sum (Y_i - \hat{Y}_i)^2 = \sum (Y_i - \hat{\alpha} - \hat{\beta}X_i)^2$$

- By partial diff of the SSE wrt $\hat{\alpha}$ and $\hat{\beta}$ and equating to zero

$$\frac{\partial SSE}{\partial \hat{\alpha}} = -2 \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i) = 0$$

$$\frac{\partial SSE}{\partial \hat{\beta}} = -2 \sum X_i (Y_i - \hat{\alpha} - \hat{\beta} X_i) = 0$$

- normal equations:

$$\sum Y_i = n\hat{\alpha} + \hat{\beta} \sum X_i$$

$$\sum X_i Y_i = \hat{\alpha} \sum X_i + \hat{\beta} \sum X_i^2$$

- Thus, 2 equations with 2 unknowns $\hat{\alpha}$ and $\hat{\beta}$. Solving for $\hat{\alpha}$ and $\hat{\beta}$

$$\hat{\beta} = \frac{n \sum X_i Y_i - (\sum X_i)(\sum Y_i)}{n \sum X_i^2 - (\sum X_i)^2} = \frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{\sum X_i^2 - n \bar{X}^2}$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}$$

$$\hat{\beta} = \frac{n \sum X_i Y_i - (\sum X_i)(\sum Y_i)}{n \sum X_i^2 - (\sum X_i)^2} = \frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{\sum X_i^2 - n \bar{X}^2}$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}$$

where $\bar{X} = \frac{1}{n} \sum X_i$ and $\bar{Y} = \frac{1}{n} \sum Y_i$

- Hence, $\hat{\alpha}$ and $\hat{\beta}$ are **ordinary least-squares (OLS) estimators** of α and β .

- The line $\hat{Y} = \hat{\alpha} + \hat{\beta} X$ is called the **least squares line** or the **estimated regression line** of Y on X.

Model in deviations form

$$Y_i = \alpha + \beta X_i + u_i \quad (2.10)$$

Applying summation & dividing by n to both sides:

$$\begin{aligned}\sum_{i=1}^n \frac{Y_i}{n} &= \sum_{i=1}^n \frac{\alpha}{n} + \sum_{i=1}^n \frac{\beta X_i}{n} + \sum_{i=1}^n \frac{u_i}{n} \\ \Rightarrow \bar{Y} &= \alpha + \beta \bar{X} + \bar{u}\end{aligned}\tag{2.11}$$

Subtracting equation (2.11) from (2.10) we get:

$$Y_i - \bar{Y} = \beta(X_i - \bar{X}) + (u_i - \bar{u})\tag{2.12}$$

Letting $x_i = X_i - \bar{X}$, $y_i = Y_i - \bar{Y}$ and $\mathcal{E}_i = (u_i - \bar{u})$, equation (2.12) becomes:

$$y_i = \beta x_i + \mathcal{E}_i\tag{2.13}$$

Equation (2.13) is the simple linear regression model in **deviations form**.

- The OLS estimator of β from equation (2.13) is given by:

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

The Gauss-Markov Theorem

- Under assumptions (1) – (5) of the **CLRM**, the OLS estimators $\hat{\alpha}$ and $\hat{\beta}$ are **Best Linear Unbiased Estimators (BLUE)**.
- The theorem tells us that of all estimators of α and β **which are linear and which are unbiased**, the estimators resulting from OLS have the **minimum variance**, that is, $\hat{\alpha}$ and $\hat{\beta}$ are the best (most efficient) linear unbiased estimators (BLUE) of α and β .

Note: If some of the assumptions stated above do not hold, then OLS estimators are **no more BLUE!!!**

Here we will prove that $\hat{\beta}$ is the BLUE of β . The proof for $\hat{\alpha}$ can be done similarly

a) To show that $\hat{\beta}$ is a linear estimator

The OLS estimator of β can be expressed as:

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} = \sum a_i y_i$$

where $a_i = \frac{x_i}{\sum x_i^2}$, $x_i = X_i - \bar{X}$ and $y_i = Y_i - \bar{Y}$.

- Thus, we can see that $\hat{\beta}$ is a **linear estimator** as it can be written as a **weighted average** of the individual observations on Y.

b) To show that $\hat{\beta}$ is an unbiased estimator of β

Note: An estimator $\hat{\theta}$ of θ is said to be **unbiased** if: $E(\hat{\theta}) = \theta$.

Consider the model in deviations form: $y_i = \beta x_i + \mathcal{E}_i$.

$$\begin{aligned}\hat{\beta} &= \frac{\sum x_i y_i}{\sum x_i^2} \\ &= \frac{\sum x_i (\beta x_i + \mathcal{E}_i)}{\sum x_i^2} \\ &= \frac{\beta \sum x_i^2 + \sum x_i \mathcal{E}_i}{\sum x_i^2} \\ &= \beta + \frac{\sum x_i \mathcal{E}_i}{\sum x_i^2} \quad (*)\end{aligned}$$

Now we have:

- $E(\beta) = \beta$ (since β is a constant)

- $E(\sum x_i \mathcal{E}_i) = \sum x_i E(\mathcal{E}_i) = \sum x_i (0) = 0$ (since x_i is non-stochastic (A5), and $E(\mathcal{E}_i) = 0$ (A2))

Thus from (*):

$$E(\hat{\beta}) = \beta + E\left[\frac{\sum x_i \mathcal{E}_i}{\sum x_i^2}\right] = \beta + \frac{\sum x_i E(\mathcal{E}_i)}{\sum x_i^2} = \beta + 0 = \beta$$

$\Rightarrow \hat{\beta}$ is an unbiased estimator of β .

$\Rightarrow \hat{\beta}$ is an unbiased estimator of β .

c) To show that $\hat{\beta}$ has the smallest variance out of all linear unbiased estimators of β

Note:

1. The OLS estimators $\hat{\alpha}$ and $\hat{\beta}$ are calculated from a **specific sample** of observations of the dependent and independent variables. If we consider a **different sample** of observations for Y and X, we get different values for $\hat{\alpha}$ and $\hat{\beta}$. This means that the values of $\hat{\alpha}$ and $\hat{\beta}$ **may vary from one sample to another**, and hence, are **random variables**.

2. The **variance** an estimator (a random variable) $\hat{\theta}$ of θ is given by:

$$Var(\hat{\theta}) = E(\hat{\theta} - \theta)^2$$

3. The expression: $(\sum_{i=1}^n x_i)^2$ can be written in expanded form as:

$$\left(\sum_{i=1}^n x_i\right)^2 = \sum_{i=1}^n x_i^2 + \sum_{i \neq j}^n x_i x_j$$

This is simply the sum of squares (x_i^2) plus the sum of cross-product terms ($x_i x_j$ for $i \neq j$).

From equation (*) we have:

$$\hat{\beta} - \beta = \frac{\sum x_i \mathcal{E}_i}{\sum x_i^2}$$

The variance of $\hat{\beta}$ is thus given by:

$$\begin{aligned} Var(\hat{\beta}) &= E(\hat{\beta} - \beta)^2 = E\left(\frac{\sum x_i \mathcal{E}_i}{\sum x_i^2}\right)^2 \\ &= \frac{1}{(\sum x_i^2)^2} E\left(\sum_{i=1}^n x_i^2 \mathcal{E}_i^2 + \sum_{i \neq j}^n x_i \mathcal{E}_i x_j \mathcal{E}_j\right) \\ &= \frac{1}{(\sum x_i^2)^2} \left(\sum_{i=1}^n x_i^2 E(\mathcal{E}_i^2) + \sum_{i \neq j}^n x_i x_j E(\mathcal{E}_i \mathcal{E}_j)\right) \\ &= \frac{1}{(\sum x_i^2)^2} \left(\sum_{i=1}^n x_i^2 (\sigma^2) + \sum_{i \neq j}^n x_i x_j (0)\right) \quad \dots\dots(**) \\ &= \frac{1}{(\sum x_i^2)^2} \left(\sigma^2 \sum_{i=1}^n x_i^2\right) = \frac{\sigma^2}{\sum x_i^2} \end{aligned}$$

Note that (**) follows from A3 and A4, that is, $var(\mathcal{E}_i) = E(\mathcal{E}_i^2) = \sigma^2$ for all i and

$cov(\mathcal{E}_i, \mathcal{E}_j) = (\mathcal{E}_i \mathcal{E}_j) = 0$ for $i \neq j$.

Thus,

$$Var(\hat{\beta}) = \frac{\sigma^2}{\sum X_i^2}$$

We have seen above (in proof (a)) that the OLS estimator of β can be expressed as:

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} = \sum a_i y_i$$

where $a_i = \frac{x_i}{\sum x_i^2}$.

Now let β^* be another linear unbiased estimator of β given by:

$$\beta^* = \sum c_i y_i$$

where $c_i = \frac{x_i}{\sum x_i^2} + d_i$ and d_i are arbitrary constants (real numbers).

β^* can be written as:

$$\begin{aligned} \beta^* &= \sum c_i y_i = \sum \left(\frac{x_i}{\sum x_i^2} + d_i \right) (\beta x_i + \mathcal{E}_i) \quad (\text{since } y_i = \beta x_i + \mathcal{E}_i) \\ &= \beta \frac{\sum x_i^2}{\sum x_i^2} + \beta \sum d_i x_i + \frac{\sum x_i \mathcal{E}_i}{\sum x_i^2} + \sum d_i \mathcal{E}_i \end{aligned}$$

Taking expectations we have:

$$\begin{aligned} E(\beta^*) &= E \left(\beta + \beta \sum d_i x_i + \frac{\sum x_i \mathcal{E}_i}{\sum x_i^2} + \sum d_i \mathcal{E}_i \right) \\ &= \beta + \beta \sum d_i x_i \quad (\text{since } E(x_i \mathcal{E}_i) = x_i E(\mathcal{E}_i) = 0, \text{ and } E(d_i \mathcal{E}_i) = d_i E(\mathcal{E}_i) = 0) \end{aligned}$$

Thus, for β^* to be unbiased (that is, for $E(\beta^*) = \beta$ to hold) we should have:

$$\sum d_i x_i = 0 \quad \dots\dots\dots(*) * *)$$

The variance of β^* is given by:

$$\begin{aligned}
Var(\beta^*) &= E(\beta^* - \beta)^2 = E(\sum c_i \mathcal{E}_i)^2 \\
&= E\left(\sum c_i^2 \mathcal{E}_i^2 + \sum_{i \neq j} c_i c_j \mathcal{E}_i \mathcal{E}_j\right) = \sum c_i^2 E(\mathcal{E}_i^2) + \sum_{i \neq j} c_i c_j E(\mathcal{E}_i \mathcal{E}_j) \\
&= \sum c_i^2 \sigma^2 + \sum_{i \neq j} c_i c_j (0) = \sigma^2 \sum c_i^2 \\
&= \sigma^2 \sum \left(\frac{x_i}{\sum x_i^2} + d_i\right)^2 = \sigma^2 \sum \left(\frac{x_i^2}{(\sum x_i^2)^2} + \frac{2x_i d_i}{\sum x_i^2} + d_i^2\right) \\
&= \sigma^2 \frac{\sum x_i^2}{(\sum x_i^2)^2} + \sigma^2 \frac{2 \sum x_i d_i}{\sum x_i^2} + \sigma^2 \sum d_i^2 \quad (\text{but } \sum d_i x_i = 0 \text{ from } (**)) \\
&= \frac{\sigma^2}{\sum x_i^2} + \sigma^2 \sum d_i^2 = Var(\hat{\beta}) + \sigma^2 \sum d_i^2
\end{aligned}$$

Thus, we have shown that:

$$Var(\beta^*) = Var(\hat{\beta}) + \sigma^2 \sum d_i^2$$

Since $\sum d_i^2$ (which is a sum of squares of real numbers) is always greater than or equal to zero, we have:

$$Var(\beta^*) \geq Var(\hat{\beta})$$

This implies that the **variance of $\hat{\beta}$ is the smallest** as compared to the variance of any other linear unbiased estimator of β .

Hence, we conclude that **$\hat{\beta}$ is the BLUE of β**

2.1.2.2 II. Maximum likelihood (ML) method of estimation Probability distribution of error terms

The OLS estimators $\hat{\alpha}$ and $\hat{\beta}$ are both linear functions of the error term, which is random by

assumption. For example:

$$\hat{\beta} = \beta + \frac{\sum x_i \mathcal{E}_i}{\sum x_i^2} = \beta + \sum a_j \mathcal{E}_i$$

where $a_i = \frac{x_i}{\sum x_i^2}$ and $x_i = X_i - \bar{X}$.

Therefore, the probability distributions of the OLS estimators will depend upon the assumptions made about the probability distribution of the error term. The nature of the probability distribution of the error term is important for hypothesis testing (or for making inferences about α and β) and also for estimation purposes.

In regression analysis, it is usually assumed that the error terms follow the normal distribution with mean 0 and variance σ^2 .

Since $\mathcal{E}_i = Y_i - \alpha - \beta X_i$, the probability distribution of \mathcal{E}_i would be:

$$\begin{aligned} p(\mathcal{E}_i) &= \frac{1}{\sqrt{2\pi}\sigma^2} \exp \left\{ -\frac{1}{2} \left(\frac{\mathcal{E}_i - E(\mathcal{E}_i)}{sd(\mathcal{E}_i)} \right)^2 \right\} \\ &= \frac{1}{\sqrt{2\pi}\sigma^2} \exp \left\{ -\frac{1}{2} \left(\frac{\mathcal{E}_i - (0)}{sd(\sigma)} \right)^2 \right\} \\ &= \frac{1}{\sqrt{2\pi}\sigma^2} \exp \left\{ -\frac{1}{2\sigma^2} (Y_i - \alpha - \beta X_i)^2 \right\} \end{aligned} \quad (2.14)$$

Here $sd(\mathcal{E}_i)$ is the standard deviation of \mathcal{E}_i , that is, $sd(\mathcal{E}_i) = \sqrt{Var(\mathcal{E}_i)} = \sigma$

Consider the linear model: $Y_i = \alpha + \beta X_i + \mathcal{E}_i$. Under the assumption that the error terms \mathcal{E}_i follow the normal distribution with mean 0 and variance σ^2 , Y_i is also normally distributed with:

$$\begin{aligned} Mean &= E(Y_i) = E(\alpha + \beta X_i + \mathcal{E}_i) = \alpha + \beta X_i \\ Variance &= Var(Y_i) = Var(\alpha + \beta X_i + \mathcal{E}_i) = Var(\mathcal{E}_i) = \sigma^2 \end{aligned}$$

Thus, the probability distribution of Y_i can be written as:

$$\begin{aligned} p(Y_i) &= \frac{1}{\sqrt{2\pi}\sigma^2} \exp \left\{ -\frac{1}{2} \left(\frac{Y_i - E(Y_i)}{sd(Y_i)} \right)^2 \right\} \\ &= \frac{1}{\sqrt{2\pi}\sigma^2} \exp \left\{ -\frac{1}{2\sigma^2} (Y_i - \alpha - \beta X_i)^2 \right\} \end{aligned}$$

ML estimation focuses on the fact that different populations generate different samples, and any one sample being scrutinized is more likely to have come from some population than from others. The ML estimator of a parameter β is the value of $\hat{\beta}$ which would most likely generate the observed sample observations Y_1, Y_2, \dots, Y_n . The ML estimator maximizes the likelihood function L which is the product of the individual probabilities (since Y_1, Y_2, \dots, Y_n are randomly selected implying independence) taken over all n observations given by:

$$\begin{aligned} L(Y_1, Y_2, \dots, Y_n, \alpha, \beta, \sigma^2) &= P(Y_1)P(Y_2)\dots P(Y_n) \\ &= \frac{1}{(\sqrt{2\pi}\sigma^2)^n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \alpha - \beta X_i)^2 \right\} \end{aligned}$$

Our aim is to maximize this likelihood function L with respect to the parameters α, β and σ^2 . To do this, it is more convenient to work with the natural logarithm of L (called **the log-likelihood function**) given by:

$$\log L = -\frac{n}{2} \log(\sigma^2) - \frac{n}{2} \log(2\pi) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \alpha - \beta X_i)^2$$

Taking partial derivatives of $\log L$ with respect to α, β & σ^2 and equating to zero, we get the ML estimators.

By partial differentiation of the $\log L$ with respect to α and β and equating the results to zero we get:

$$\begin{aligned} \frac{\partial \log L}{\partial \alpha} &= -\frac{1}{2\sigma^2} \sum (Y_i - \alpha - \beta X_i) (-X_i) = 0 \\ \frac{\partial \log L}{\partial \beta} &= -\frac{1}{2\sigma^2} \sum (Y_i - \alpha - \beta X_i) (-1) = 0 \end{aligned}$$

Re-arranging the two equations, and replacing β by β_{ML} and α by α_{ML} , we get:

$$\beta_{ML} = \frac{n \sum X_i Y_i - (\sum X_i)(\sum Y_i)}{n \sum X_i^2 - (\sum X_i)^2} = \frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{\sum X_i^2 - n \bar{X}^2} = \hat{\beta}$$

$$\alpha_{ML} = \bar{Y} - \beta_{ML} \bar{X} = \hat{\alpha}$$

By partial differentiation of the log L with respect to σ^2 and equating to zero we get:

$$\frac{\partial \log L}{\partial \sigma^2} = -\frac{n}{2} \left(\frac{1}{\sigma^2} \right) - \frac{1}{2} \sum (Y_i - \alpha - \beta X_i) \left(\frac{-1}{(\sigma^2)^2} \right) = 0$$

Replacing σ^2 by σ_{ML}^2 and simplifying, we get:

$$\sigma_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \alpha_{ML} - \beta_{ML} X_i)^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta} X_i)^2 = \frac{1}{n} \sum_{i=1}^n \hat{\mathcal{E}}_i^2$$

Note

- 1) The ML estimators $\hat{\alpha}_{ML}$ and $\hat{\beta}_{ML}$ are identical to the OLS estimators, and are thus best linear unbiased estimators (BLUE) of α and β , respectively.
- 2) The ML estimator $\hat{\sigma}_{ML}^2$ of σ^2 is biased.

Proof

$$\begin{aligned} \hat{\mathcal{E}}_i &= Y_i - \hat{Y}_i = Y_i - (\hat{\alpha} + \hat{\beta} X_i) \\ &= (\alpha + \beta X_i + \mathcal{E}_i) - (\hat{\alpha} + \hat{\beta} X_i) \\ &= (\alpha - \hat{\alpha}) + (\beta - \hat{\beta}) X_i + \mathcal{E}_i \end{aligned} \tag{2.15}$$

$$\begin{aligned} Y_i &= \alpha + \beta X_i + \mathcal{E}_i \\ \Rightarrow \bar{Y} &= \alpha + \beta \bar{X} + \bar{\mathcal{E}} \\ \Rightarrow \alpha &= \bar{Y} - \beta \bar{X} - \bar{\mathcal{E}} \end{aligned} \tag{2.16}$$

We know that the OLS estimator of α is given by:

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} \tag{2.17}$$

Subtracting (2.17) from (2.16) we get:

$$(\alpha - \hat{\alpha}) = (\hat{\beta} - \beta)\bar{X} - \bar{\mathcal{E}} = -(\beta - \hat{\beta})\bar{X} - \bar{\mathcal{E}} \quad (2.18)$$

Substituting equation (2.18) in (2.15) we get:

$$\begin{aligned} \hat{\mathcal{E}}_i &= -(\beta - \hat{\beta})\bar{X} - \bar{\mathcal{E}} + (\beta - \hat{\beta})X_i - \mathcal{E}_i = (\beta - \hat{\beta})(X_i - \bar{X}) + (\mathcal{E}_i - \bar{\mathcal{E}}) \\ &= (\beta - \hat{\beta})x_i + e_i \end{aligned}$$

where $x_i = X_i - \bar{X}$ and $e_i = \mathcal{E}_i - \bar{\mathcal{E}}$. Squaring both sides and taking summations we have:

$$\sum_{i=1}^n \hat{\mathcal{E}}_i^2 = (\beta - \hat{\beta})^2 \sum_{i=1}^n x_i^2 + 2(\beta - \hat{\beta}) \sum_{i=1}^n x_i e_i + \sum_{i=1}^n e_i^2 \quad (2.19)$$

From the two-variable model in deviations form we have:

$$\begin{aligned} \hat{\beta} - \beta &= \frac{\sum x_i \mathcal{E}_i}{\sum x_i^2} \\ &= \frac{\sum x_i (\mathcal{E}_i - \bar{\mathcal{E}})}{\sum x_i^2} \\ &= \frac{\sum x_i e_i}{\sum x_i^2} \\ (\text{since } \frac{\sum x_i \bar{\mathcal{E}}}{\sum x_i^2} &= \frac{\bar{\mathcal{E}} \sum x_i}{\sum x_i^2} = \frac{\bar{\mathcal{E}}(0)}{\sum x_i^2} = 0) \\ \Rightarrow \sum x_i e_i &= (\hat{\beta} - \beta) \sum x_i^2 = -(\beta - \hat{\beta}) \sum x_i^2 \end{aligned} \quad (2.20)$$

Substituting (2.20) in (2.19) we have:

$$\begin{aligned} \sum_{i=1}^n \hat{\mathcal{E}}_i^2 &= (\beta - \hat{\beta})^2 \sum_{i=1}^n x_i^2 - 2(\beta - \hat{\beta})^2 \sum_{i=1}^n x_i^2 + \sum_{i=1}^n e_i^2 = -(\beta - \hat{\beta})^2 \sum_{i=1}^n x_i^2 + \sum_{i=1}^n e_i^2 \\ \bullet E \left(-(\beta - \hat{\beta})^2 \sum_{i=1}^n x_i^2 \right) &= - \sum_{i=1}^n x_i^2 E \left[(\beta - \hat{\beta})^2 \right] = - \sum_{i=1}^n x_i^2 \text{Var}(\hat{\beta}) = - \sum_{i=1}^n x_i^2 \left[\frac{\sigma^2}{\sum_{i=1}^n x_i^2} \right] = -\sigma^2 \\ \bullet E \left(\sum_{i=1}^n e_i^2 \right) &= E \left(\sum_{i=1}^2 (\mathcal{E}_i - \bar{\mathcal{E}})^2 \right) = E \left(\sum_{i=1}^2 \mathcal{E}_i^2 - n\bar{\mathcal{E}}^2 \right) \\ &= \sum_{i=1}^2 E(\mathcal{E}_i^2) - nE(\bar{\mathcal{E}}^2) = n\sigma^2 - n \left(\frac{\sigma^2}{n} \right) = n\sigma^2 - \sigma^2 = (n-1)\sigma^2 \end{aligned}$$

Then it follows that:

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

Thus,

$$E(\sigma_{ML}^2) = E \left(\frac{1}{n} \sum_{i=1}^n \hat{\mathcal{E}}_i^2 \right) = \frac{1}{n} E \left(\sum_{i=1}^n \hat{\mathcal{E}}_i^2 \right) = \left(\frac{n-2}{n} \right) \sigma^2 \neq \sigma^2$$

From the above result it follows that an unbiased estimator of σ^2 is:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\mathcal{E}}_i^2$$

2.1.3 Statistical inference in simple linear regression model

Estimation of standard error

To make statistical inferences about the true (population) regression coefficient β , we make use of the estimator $\hat{\beta}$ and its variance $Var(\hat{\beta})$. We have already seen that:

$$Var(\hat{\beta}) = \frac{\sigma^2}{\sum x_i^2}$$

where $x_i = X_i - \bar{X}$. Since this variance depends on the unknown parameter σ^2 , we have to estimate σ^2 . As shown above, an unbiased estimator of σ^2 is given by:

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}X_i)^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\mathcal{E}}_i^2$$

Thus, an unbiased estimator of $Var(\hat{\beta})$ is given by:

$$\hat{Var}(\hat{\beta}) = \frac{\hat{\sigma}^2}{\sum x_i^2} = \frac{\sum \hat{\mathcal{E}}_i^2}{(n-2) \sum x_i^2}$$

The square root of $\widehat{Var}(\hat{\beta})$ is called the standard error of $\hat{\beta}$, that is,

$$s.e.(\hat{\beta}) = \sqrt{\widehat{Var}(\hat{\beta})} = \sqrt{\frac{\hat{\sigma}^2}{\sum x_i^2}}$$

Tests of significance of regression coefficients

Consider the simple linear regression model:

$$Y_i = \alpha + \beta X_i + \mathcal{E}_i$$

If there is no relationship between X and Y, then this is equivalent to saying $\beta = 0$ (β is not significantly different from zero). Thus, the null hypothesis of no relationship between X and Y is expressed as:

$$H: \beta = 0$$

The alternative hypothesis is that there is a significant relationship between X and Y, that is,

$$H: \beta \neq 0$$

In order to reject or not reject the null hypothesis, we calculate the **test statistic** given by:

$$t = \frac{\hat{\beta} - \beta_0}{s.e.(\hat{\beta})} = \frac{\hat{\beta} - 0}{s.e.(\hat{\beta})} = \frac{\hat{\beta}}{s.e.(\hat{\beta})}$$

and compare this figure with the value from the student's t distribution with (n-2) degrees of freedom for a given significance level α .

Decision rule: If $|t| > t_{\frac{\alpha}{2}}(n - 2)$ then we **reject the null hypothesis**, and conclude that there is a significant relationship between X and Y.

Confidence interval for β

Confidence interval provides a range of values which are likely to contain the true regression parameter. With every confidence interval, we associate a level of statistical significance (α). The confidence intervals are constructed in such a way that the probability of the interval to contain the true parameter is (1 - α). Symbolically,

$$\begin{aligned}
P[-t_{\frac{\alpha}{2}(n-2)} < t < t_{\frac{\alpha}{2}(n-2)}] &= 1 - \alpha \\
\Rightarrow P[-t_{\frac{\alpha}{2}(n-2)} < \frac{\hat{\beta} - \beta_0}{s.e(\hat{\beta})} < t_{\frac{\alpha}{2}(n-2)}] \\
\Rightarrow P[\hat{\beta} - t_{\frac{\alpha}{2}(n-2)} * s.e(\hat{\beta}) < \beta_0 < \hat{\beta} + t_{\frac{\alpha}{2}(n-2)} * s.e(\hat{\beta})] &= 1 - \alpha
\end{aligned}$$

Thus, a $(1 - \alpha)100\%$ confidence interval for β is given by:

$$\hat{\beta} \pm t_{\frac{\alpha}{2}(n-2)} * s.e(\hat{\beta})$$

Test of model adequacy

Is the estimated equation a useful one? To answer this, an objective measure of some sort is desirable.

The **total variation** in the dependent variable Y is given by:

$$Variation(Y) = \sum (Y_i - \bar{Y})^2$$

Our goal is to partition this variation into two: one that accounts for variation due to the regression equation (explained portion) and another that is associated with the unexplained portion of the model.

We can write $Y_i - \bar{Y}$ as:

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

Squaring both sides and taking summations we have:

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

Remark 1:

$$\sum_{i=1}^n \hat{\mathcal{E}}_i = 0 \text{ where } \hat{\mathcal{E}}_i = Y_i - \hat{Y}_i = Y_i - \hat{\alpha} - \hat{\beta}X_i$$

Proof:

$$\begin{aligned} \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}X_i) &= \sum Y_i - \sum \hat{\alpha} - \hat{\beta} \sum Y_i \\ &= n\bar{Y} - n\hat{\alpha} - n\hat{\beta}\bar{X} \\ &= n[(\bar{Y} - \hat{\beta}) - \hat{\alpha}] = n(\hat{\alpha} - \hat{\alpha}) \end{aligned}$$

Remark 2:

$$\sum_{i=1}^n \hat{\mathcal{E}}_i X_i = 0$$

Proof:

$$\begin{aligned} \sum_{i=1}^n \hat{\mathcal{E}}_i X_i &= \sum (Y_i - \hat{\alpha} - \hat{\beta}X_i) X_i \\ &= \sum Y_i X_i - \sum \hat{\alpha} X_i - \sum \hat{\beta} X_i^2 \\ &= \sum Y_i X_i - \sum \hat{\alpha} X_i - \sum \hat{\beta} X_i^2 \\ &= \sum Y_i X_i - n\hat{\alpha}\bar{X} - \sum \hat{\beta} X_i^2 \\ &= \sum Y_i X_i - n(\bar{Y} - \hat{\beta}\bar{X})\bar{X} - \sum \hat{\beta} X_i^2 \\ &= \sum Y_i X_i - n\bar{X}\bar{Y} - \hat{\beta}[\sum X_i^2 - n\bar{X}^2] \\ &= \sum Y_i X_i - n\bar{X}\bar{Y} - \frac{\sum Y_i X_i - n\bar{X}\bar{Y}}{\sum X_i^2 - n\bar{X}^2} [\sum X_i^2 - n\bar{X}^2] = 0 \end{aligned}$$

From remarks 1 and 2 it follows that:

$$\begin{aligned} \sum (Y_i - \hat{Y}_i)(\hat{Y}_i - \bar{Y}) &= \sum \hat{\mathcal{E}}_i \hat{Y}_i - \sum \hat{\mathcal{E}}_i \bar{Y} \\ &= \sum \hat{\mathcal{E}}_i (\hat{\alpha} + \hat{\beta}X_i) - \bar{Y} \underbrace{\sum \hat{\mathcal{E}}_i}_{=0} = \hat{\alpha} \underbrace{\sum \hat{\mathcal{E}}_i}_{=0} + \hat{\beta} \underbrace{\sum \hat{\mathcal{E}}_i X_i}_{=0} = 0 \end{aligned}$$

Thus, the cross product term in equation (*) vanishes, and we are left with:

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

<i>Variation in Y</i>	<i>Residual variation</i>	<i>Explained variation</i>
TSS	$=$	ESS
		$+ \quad RSS$

In other words, the total sum of squares (TSS) is decomposed into regression (explained) sum of squares (RSS) and error (residual or unexplained) sum of squares (ESS).

Computational formulas:

•The total sum of squares (TSS) is a measure of dispersion of the observed values of Y about their mean. This is computed as:

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

•The regression (explained) sum of squares (RSS) measures the amount of the total variability in the observed values of Y that is accounted for by the linear relationship between the observed values of X and Y. This is computed as:

$$RSS = \sum (\hat{Y}_i - \bar{Y})^2 = \hat{\beta}^2 \left[\sum_{i=1}^n (X_i - \bar{X})^2 \right] = \hat{\beta}^2 \sum_{i=1}^n x_i^2$$

•The error (residual or unexplained) sum of squares (ESS) is a measure of the dispersion of the observed values of Y about the regression line. This is computed as:

$$ESS = \sum (Y_i - \hat{Y}_i)^2 = TSS - RSS$$

If a regression equation does a good job of describing the relationship between two variables, the explained sum of squares should constitute a large proportion of the total sum of squares. Thus, it would be of interest to determine the magnitude of this proportion by computing the ratio of the explained sum of squares to the total sum of squares. This proportion is called the sample

coefficient of determination R^2 . That is:

$$\text{Coefficient of determination} = R^2 = \frac{RSS}{TSS} = 1 - \frac{ESS}{TSS}$$

The coefficient of determination can also be computed as:

$$R^2 = \frac{\hat{\beta} \sum x_i y_i}{\sum y_i^2} \text{ where } x_i = X_i - \bar{X} \text{ and } y_i = Y_i - \bar{Y}$$

Tests for the coefficient of determination (R^2)

The largest value that R^2 can assume is 1 (in which case all observations fall on the regression line), and the smallest it can assume is zero. A low value of R^2 is an indication that:

- X is a poor explanatory variable in the sense that variation in X leaves Y unaffected, or
- while X is a relevant variable, its influence on Y is weak as compared to some other variables that are omitted from the regression equation, or
- the regression equation is misspecified (for example, an exponential relationship might be more appropriate).

Note

1) The proportion of total variation in the dependent variable (Y) that is explained by changes in the independent variable (X) or by the regression line is equal to: $R^2 * 100\%$.

2) The proportion of total variation in the dependent variable (Y) that is due to factors other than X (for example, due to excluded variables, chance, etc) is equal to: $(1 - R^2)100\%$

Thus, a small value of R^2 casts doubt about the usefulness of the regression equation. We do not, however, pass final judgment on the equation until it has been subjected to an objective statistical test. Such a test is accomplished by means of analysis of variance (ANOVA) which enables us to test the significance of R^2 (i.e., the adequacy of the linear regression model).

The ANOVA table for simple linear regression is given below:

ANOVA table

Source of variation	Sum of squares	Degrees of freedom	Mean square	Variance ratio
Regression	RSS	1	RSS/1	$F_{cal} = \frac{RSS/1}{ESS/(n-2)}$
Residual	ESS	n-2	$\frac{ESS}{(n-2)}$	
Total	TSS	n-1		

To test for the significance of R^2 , we compare the variance ratio with the critical value from the F

distribution with 1 and (n-2) degrees of freedom in the numerator and denominator, respectively, for a given significance level α . **Decision:** If the calculated variance ratio exceeds the tabulated value, that is, if $F_{cal} > F_{\alpha(1,n-2)}$, we then conclude that R^2 is significant (or that the linear regression model is adequate).

Note: The F test is designed to test the significance of all variables or a set of variables in a regression model. In the two-variable model, however, it is used to test the explanatory power of a single variable (X), and at the same time, is equivalent to the test of significance of R^2 .

Example

Consider the following data on the percentage rate of change in electricity consumption (millions KWH) (Y) and the rate of change in the price of electricity (Birr/KWH) (X) for the years 1979 – 1994. Use **chp1.csv**.

```
#electricity consumption
ele_cons=read.csv("ele_cons.csv",header = T,sep = ",")
(n<-length(ele_cons))

## [1] 3

x_i<-ele_cons$X-mean(ele_cons$X)
y_i<-ele_cons$Y-mean(ele_cons$Y)
(x_bar<-mean(ele_cons$X))

## [1] 1.281

(y_bar<-mean(ele_cons$Y))

## [1] 23.43

(sum_x_i_sq<-sum(x_i^2)) # sum of square of deviation of x

## [1] 92.2

(sum_y_i_sq<-sum(y_i^2)) # sum of square of deviation of y

## [1] 13229
```

```
(sum_xy<-sum(x_i*y_i)) #Sum of cross product of x*y
```

```
## [1] -779.2
```

```
(b_hat<-sum_xy/sum_x_i_sq) # Slope using deviation form
```

```
## [1] -8.451
```

```
(a_hat<-y_bar-b_hat*x_bar)
```

```
## [1] 34.25
```

Therefore, the estimated regression equation is:

$$\hat{Y} = \hat{\alpha} + \hat{\beta}X \Leftrightarrow \hat{Y} = 34.25 + -8.451X \quad \textbf{Test of model adequacy:}$$

$$TSS = \sum y_i^2$$

```
tss=sum(y_i^2)
```

$$TSS = \sum y_i^2 = 1.3229 \times 10^4$$

$$RSS = \hat{\beta}^2 \sum x_i^2$$

```
(rss<-b_hat^2*sum(x_i^2)) # Regression Sum of Square
```

```
## [1] 6586
```

```
(ess<-tss-rss) # Error sum of square
```

```
## [1] 6643
```

```
(r_sq<-rss/tss) # Coefficient of determination
```

```
## [1] 0.4978
```

```
# Mean square Error (EMS) = ESS / (n-k)
ems <- ess / (nrow(ele_cons) - 2)

# Mean square of Regression (RMS) = RSS / (k-1)
rms <- rss / 1
(f_cal = rms / ems)

## [1] 13.88
```

$$RSS = \hat{\beta}^2 \sum x_i^2 = 6585.68$$

$$\Rightarrow R^2 = \frac{RSS}{TSS} = 6585.68 / 1.3229 \times 10^4 = 0.4978$$

Thus, we can conclude that:

About 50% of the variation in electricity consumption is due to changes in the price of electricity. The remaining 50% of the variation in electricity consumption is not due to changes in the price of electricity, but instead due to chance and other factors not included in the model.

ANOVA table

Source of variation	Sum of squares	Degrees of freedom	Mean square	Variance ratio
Regression	6585.68	1	6585.68	13.8792
Residual	6643.02	16-2=14	474.5	
Total	1.3229×10^4	16-1=15		

For $\alpha = 0.05$, the critical value from the F-distribution is:

$$F_{\alpha(1, n-2)} = F_{0.05(1, 14)}$$

```
(qf(p=0.05, df1=1, df2=14, lower.tail = F)) # Tabulated F-value

## [1] 4.6
```

```
#Or p-value
(pf(f_cal,df1=1,df2=14,lower.tail = F)) # T

## [1] 0.00226
```

Decision: Since the calculated variance ratio exceeds the critical value, we reject the null hypothesis of no linear relationship between price and consumption of electricity at the 5% level of significance. Thus, we then conclude that R^2 is significant, that is, the linear regression model is adequate and is useful for prediction purposes.

Estimation of the standard error of the coefficients and test of its significance

An unbiased estimator of the error variance

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^{16} \hat{\mathcal{E}}_i^2 = \frac{1}{n-2} \sum_{i=1}^{16} (Y_i - \hat{Y})^2 = \frac{ESS}{n-2}$$

```
y_hat<-a_hat+b_hat*ele_cons$X # Fitted values
e_hat<-ele_cons$Y-y_hat # Residuals
(sigma_hat_sq<-sum(e_hat^2)/(nrow(ele_cons)-2))

## [1] 474.5
```

```
# Or
(sigma_hat_sq<-ess/(nrow(ele_cons)-2))

## [1] 474.5
```

Thus, an unbiased estimator of $Var(\hat{\beta})$ is given by

$$Var(\hat{\beta}) = \frac{\hat{\sigma}^2}{\sum x_i^2}$$

```
(var_b_hat<-sigma_hat_sq/sum_x_i_sq)

## [1] 5.146
```

$$Var(\hat{\beta}) = \frac{\hat{\sigma}^2}{\sum x_i^2} = 474.5012/92.2011 = 5.1464$$

```
s.e_b_hat<-sqrt(var_b_hat)
```

The standard error of $\hat{\beta}$ is:

$$s.e(\hat{\beta}) = \sqrt{Var(\hat{\beta})} = \sqrt{5.146} = 2.2686$$

The hypothesis of interest is: $H_0 : \beta = 0$ VS $H_1 : \beta \neq 0$ We calculate the test statistic:

```
(t_cal<-(b_hat-0)/s.e_b_hat)
```

```
## [1] -3.725
```

For $\alpha = 0.05$, the critical value from the student's *tdistribution* with (n-2) degrees of freedom is:

```
(qt(p=0.05/2,df=nrow(ele_cons)-2,lower.tail = F))
```

```
## [1] 2.145
```

```
#Or using p-value
```

```
(pt(t_cal,df=nrow(ele_cons)-2,lower.tail = T))
```

```
## [1] 0.00113
```

Decision: Since $|t| > t_{\alpha \frac{1}{2}(n-2)}$, we reject the null hypothesis, and conclude that β is significantly different from zero. In other words, the price of electricity significantly and negatively affects electricity consumption.

The interpretation of the estimated regression coefficient $\hat{\beta}=-8.451$ is that for a one percent drop (increase) in the growth rate of price of electricity, there is an 8.45 percent increase (decrease) in the growth rate of electricity consumption.

2.2 Multiple Linear Regressions

So far we have seen the basic statistical tools and procedures for analyzing relationships between two variables. But in practice, economic models generally contain one dependent variable and two or more independent variables. Such models are called **multiple regression models**.

Economic relationships usually include more than one regressor. For example, a demand equation for a product will usually include real price of that product(\hat{Y}_i) in addition to real income as well as real price of a competitive product and the advertising expenditures on the product. In this case the sales of the product is modelled as

$$\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3$$

where Y_i denotes the i^{th} observation on the dependent variable Y, in this case the sales of this product. own price(X_1), the competitor's price(X_2) and advertising expenditures(X_3)

Example: a) In demand studies we study the relationship between the demand for a good (Y) and price of the good (X_2), prices of substitute goods (X_3) and the consumer's income (X_4). Here, Y is the dependent variable and X_2 , X_3 and X_4 are the explanatory (independent) variables. The relationship is estimated by a multiple linear regression equation (model) of the form:

$$\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3 + \hat{\beta}_4 X_4$$

where $\hat{\beta}_1$, $\hat{\beta}_2$, $\hat{\beta}_3$ and $\hat{\beta}_4$ are estimated regression coefficients.

- b) In a study of the amount of output (product), we are interested to establish a relationship between output (Q) and labour input (L) ' & ' capital input (K). The equations are often estimated in log-linear form as:

$$\log(\hat{Q}) = \hat{\beta}_1 + \hat{\beta}_2 \log(L) + \hat{\beta}_3 \log(K)$$

- c) In a study of the determinants of the number of children born per woman (Y), the possible explanatory variables include years of schooling of the woman (X_2), woman's (or husband's)

earning at marriage (X_3), age of woman at marriage (X_4) and survival probability of children at age five (X_5). The relationship can thus be expressed as:

$$\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3 + \hat{\beta}_4 X_4 + \hat{\beta}_5 X_5$$

Model assumptions

Dependent variable: Y of size $n \times 1$

Independent (explanatory) variables: X_2, X_3, \dots, X_k each of size $n \times 1$

Assumptions

1. The true model is: $Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k X_{ki} + \mathcal{E}_i$.
2. The error terms have **zero mean**: $E(\mathcal{E}_i)$.
3. **Homoscedasticity**: $var(\mathcal{E}_i) = E(\mathcal{E}_i^2) = \sigma^2$ for all i .
4. **No error autocorrelation**: $cov(\mathcal{E}_i, \mathcal{E}_j) = E(\mathcal{E}_i \mathcal{E}_j) = 0$ for $i \neq j$.
5. Each of the explanatory variables X_2, X_3, \dots, X_k is **non-stochastic**.
6. **No multicollinearity**: No exact linear relationship exists between any of the explanatory variables.
7. **Normality**: \mathcal{E}_i are normally distributed with mean zero and variance σ^2 for all i ($\mathcal{E}_i \sim N(0, \sigma^2)$).

The only additional assumption here is that there is no multicollinearity, meaning that there is no linear dependence between the regressor variables X_2, X_3, \dots, X_k .

Under the above assumptions, ordinary least squares (OLS) yields best linear unbiased estimators (BLUE) of $\beta_1, \beta_2, \dots, \beta_k$.

2.2.1 Estimation of parameters and standard errors

Example: Consider the following model ($K = 3$):

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \mathcal{E}_i \quad (2.21)$$

$$\frac{\sum Y_i}{n} = \beta_1 + \beta_2 \frac{\sum X_{2i}}{n} + \frac{\mathcal{E}_i}{n} + \beta_3 \frac{\sum X_{3i}}{n}$$

$$\bar{y} = \beta_1 + \beta_2 \bar{x}_{2i} + \beta_3 \bar{x}_{3i} + \bar{\mathcal{E}}. \quad (2.22)$$

Equation (2.21)-(2.22)

$$Y_i - \bar{y} = \beta_1 - \beta_1 + \beta_2(X_{2i} - \bar{X}_2) + \beta_3(X_{3i} - \bar{X}_3) + (\mathcal{E}_i - \bar{\mathcal{E}})$$

Hence the model in deviation form:

$$y_i = \beta_2 x_{2i} + \beta_3 x_{3i} + e_i \text{ where } e_i = \mathcal{E}_i - \bar{\mathcal{E}}$$

OLS estimator of β_2 and β_3 is an estimator which minimizes the sum square of error(ESS) in the deviation form as follows:

$$ESS = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \beta_2 x_{2i} - \beta_3 x_{3i})^2$$

The error sum of squares (ESS) is:

$$ESS = \sum \mathcal{E}_i^2 = \sum (y_i - \beta_2 X_{2i} - \beta_3 X_{3i})^2$$

Partially differentiating the ESS with respect to β_2 and β_3 , equating to zero and simplifying, we get:

$$\begin{aligned}
\frac{\partial ESS}{\partial \beta_2} = 0 &\Rightarrow 2 \sum (y_i - \beta_2 x_{2i} - \beta_3 x_{3i})(-x_{2i}) = 0 \\
&\Rightarrow \sum y_i x_{2i} - \beta_2 \sum x_{2i}^2 - \beta_3 \sum x_{3i} x_{2i} = 0 \\
&\quad \sum y_i x_{2i} = \beta_2 \sum x_{2i}^2 + \beta_3 \sum x_{3i} x_{2i}
\end{aligned} \tag{2.23}$$

$$\begin{aligned}
\frac{\partial ESS}{\partial \beta_3} = 0 &\Rightarrow 2 \sum (y_i - \beta_2 x_{2i} - \beta_3 x_{3i})(-x_{3i}) = 0 \\
&\Rightarrow \sum y_i x_{3i} - \beta_2 \sum x_{2i} x_{3i} - \beta_3 \sum x_{3i}^2 = 0 \\
&\quad \sum y_i x_{3i} = \beta_2 \sum x_{2i} x_{3i} + \beta_3 \sum x_{3i}^2
\end{aligned} \tag{2.24}$$

Using normal equation (2.23):

$$\beta_2 = \frac{\sum y_i x_{2i} - \beta_3 \sum x_{3i} x_{2i}}{\sum x_{2i}^2} \tag{2.25}$$

Substituting value of β_2 in normal equation (2.24):

$$\begin{aligned}
\sum y_i x_{3i} &= \left(\frac{\sum y_i x_{2i} - \beta_3 \sum x_{3i} x_{2i}}{\sum x_{2i}^2} \right) \sum x_{2i} x_{3i} + \beta_3 \sum x_{3i}^2 \\
\Rightarrow \hat{\beta}_3 &= \frac{[\sum x_{3i} y_i][\sum x_{2i}^2] - [\sum x_{2i} y_i][\sum x_{2i} x_{3i}]}{[\sum x_{2i}^2][\sum x_{3i}^2] - [\sum x_{2i} x_{3i}]^2}
\end{aligned} \tag{2.26}$$

Substituting value of $\hat{\beta}_3$ in equation ((2.25))

$$\hat{\beta}_2 = \frac{[\sum x_{2i} y_i][\sum x_{3i}^2] - [\sum x_{3i} y_i][\sum x_{2i} x_{3i}]}{[\sum x_{2i}^2][\sum x_{3i}^2] - [\sum x_{2i} x_{3i}]^2}$$

An estimator of β_1 is also obtained from the non-deviation form:

$$\frac{\partial ESS}{\partial \beta_1} = 0 \implies \frac{\partial \sum \mathcal{E}_i^2}{\partial \beta_1} = 2 \sum (Y_i - \beta_1 - \beta_2 X_{2i} - \beta_3 X_{3i})(-1) = 0$$

$$\sum Y_i = n\beta_1 + \beta_2 \sum X_{2i} + \beta_3 \sum X_{3i}$$

Dividing both sides by n :

$$\implies \hat{\beta} = \bar{Y} - \hat{\beta}_2 \bar{X}_2 - \hat{\beta}_3 \bar{X}_3$$

2.2.1.1 Variances of estimated regression coefficients

- First we have to determine the **unbiased estimator of the variance of the errors** σ^2 given by:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{\mathcal{E}}_i^2}{n - k} = \frac{\sum_{i=1}^n \hat{\mathcal{E}}_i^2}{n - 3} = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n - 3}$$

where $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i}$

- $\hat{\beta}_2$, the OLS estimator of β_2 , the **partial derivative** of Y_i with respect to X_{2i}
- we can interpret $\hat{\beta}_2$ as a simple linear regression coefficient.

Steps

1. Run the regression of X2 **on all the other X's** , and obtain the residuals \hat{v}_2 , i.e.,

$$X_{2i} = \hat{X}_{2i} + \hat{v}_2$$
2. Run the simple regression of Y on \hat{v}_2 , the resulting estimate of the slope coefficient is $\hat{\beta}_2$
 - The first regression essentially **cleans out the effect of the other X's from X2**, leaving the variation **unique to X2** in \hat{v}_2
 - Then

$$var(\hat{\beta}_2) = \sigma^2 / \sum_{i=1}^n \hat{v}_{2i}^2 \quad (2.27)$$

- Let R_2^2 be the R^2 for the **regression of X_2 on all the other X 's**, then

$$R_2^2 = 1 - \frac{ESS}{TSS} = 1 - \frac{\sum_{i=1}^n \hat{v}_{2i}^2}{\sum_{i=1}^n x_{2i}^2} \Rightarrow \sum_{i=1}^n \hat{v}_{2i}^2 = \sum_{i=1}^n x_{2i}^2 (1 - R_2^2) \quad (2.28)$$

$$Var(\hat{\beta}_2) = \frac{\sigma^2}{\sum_{i=1}^n \hat{v}_{2i}^2} = \frac{\sigma^2}{\sum_{i=1}^n x_{2i}^2 (1 - R_2^2)} \quad (2.29)$$

For case of $k=3$

- Or use

$$\hat{V}(\hat{\beta}_2) = \frac{\hat{\sigma}^2}{(1 - r_{23}^2) \sum X_{2i}^2} \quad (2.30)$$

$$\text{and } \hat{V}(\hat{\beta}_3) = \frac{\hat{\sigma}^2}{(1 - r_{23}^2) \sum X_{3i}^2} \quad (2.31)$$

where r_{23} is the coefficient of correlation between X_2 and X_3 , that is:

$$r_{23} = \frac{\sum x_{2i}x_{3i}}{\sqrt{(\sum x_{2i}^2)(\sum x_{3i}^2)}}$$

Taking the square roots, we obtain the standard errors of $\hat{\beta}_2$ and $\hat{\beta}_3$: $\text{s.e.}(\hat{\beta}_2) = \sqrt{\hat{V}(\hat{\beta}_2)}$,
 $\text{s.e.}(\hat{\beta}_3) = \sqrt{\hat{V}(\hat{\beta}_3)}$.

2.2.2 The coefficient of determination and test of model adequacy

The coefficient of determination (R^2) can be calculated as usual as:

$$R^2 = \frac{RSS}{TSS} = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\sum_{i=1}^n \hat{\mathcal{E}}_i^2}{\sum_{i=1}^n Y_i^2}$$

R^2 measures the proportion of variation in the dependent variable Y that is explained by the explanatory variables (or by the multiple linear regression model). It is a **goodness-of-fit statistic**.

A test for the significance of R^2 or a test of model adequacy is accomplished by testing the hypotheses:

$$H_0 : \beta_2 = \beta_3 = 0$$

$$H_A : H_0 \text{ at least one } \beta_i \neq 0$$

The test statistic is given by:

$$F_{cal} = \frac{RSS/(k-1)}{TSS/(n-k)} = \frac{RSS/(3-1)}{TSS/(n-3)}$$

where K is the number of parameters estimated from the sample data (K = 3 in our case since we estimate β_1, β_2 and β_3) and n is the sample size. We say that the linear model is adequate in explaining the relationship between the dependent variable and one or more of the independent variables if:

$$F_{cal} = F_{\alpha}(k-1, n-k)$$

- Since **OLS minimizes** the residual sums of squares, **adding one or more variables to the regression cannot increase** this residual sums of squares
- Makes $\sum \hat{\mathcal{E}}^2$ non-increasing and R^2 non-decreasing, since $R^2 = 1 - \frac{\sum \hat{\mathcal{E}}^2}{\sum y_i^2}$
- $\bar{R}^2 = 1 - [\frac{\sum \hat{\mathcal{E}}^2/(n-k)}{\sum y_i^2/(n-1)}]$ adjusted by their degrees of freedom
- this variable will increase \bar{R}^2 **only if the reduction** in $\sum \hat{\mathcal{E}}^2$ **outweighs this loss**

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

Unlike R^2 , \bar{R}^2 may increase or decrease when new variables are added into the model.

2.2.3 Tests on the regression coefficients

To test whether each of the coefficients are significant or not, the null and alternative hypotheses are given by:

$$H_0 : \beta_j = 0$$

$$H_A : \beta_j \neq 0$$

Table 2.2: Food Consumption

Year	Y	X2	X3
1927	88.9	91.7	57.7
1928	88.9	92.0	59.3
1929	89.1	93.1	62.0
1930	88.7	90.9	56.3
1931	88.0	82.3	52.7
1932	85.9	76.3	44.4

for $j = 2, 3$. The test statistic is:

$$t_j = \frac{\hat{\beta}_j}{s.e.(\hat{\beta}_j)} \quad j = 2, 3$$

Decision rule:

If $|t_j| > t_{\alpha/2}(n - 3)$, we reject H_0 and conclude that β_j is significant, that is, the regressor variable X_j , $j=2, 3$ significantly affects the dependent variable Y .

Example

Consider the following data on per capita food consumption (Y), price of food (X_2) and per capital income (X_3) for the years 1927-1941 in the United States. Retail price of food and per capital disposable income are deflated by the Consumer Price Index and the data is stored as `food_cons.csv` (see the first six observations in Table 2.2). Fit a multiple linear regression model:

```
food_cons<-read.csv("food_cons.csv", header = T, sep=",")
```

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \mathcal{E}_i, i = 1, 2, \dots, 15$$

To simplify the calculations, it is better to work with deviations: $y_i = Y_i - \bar{Y}$, $x_{2i} = X_{2i} - \bar{X}_2$ and $x_{3i} = X_{3i} - \bar{X}_3$.

Table 2.3: In deviation form

x2i	x3i	yi
5.8	1.1733	-0.0067
6.1	2.7733	-0.0067
7.2	5.4733	0.1933
5.0	-0.2267	-0.2067
-3.6	-3.8267	-0.9067
-9.6	-12.1267	-3.0067

Table 2.4: Estiamted Coefficents

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	86.0832	3.8730	22.226	0.0000
X2	-0.2160	0.0520	-4.155	0.0013
X3	0.3781	0.0338	11.178	0.0000

```
#Writing in Deviation form(yi,x2i,x3i)
food_cons$x2i=food_cons$X2-mean(food_cons$X2)
food_cons$x3i=food_cons$X3-mean(food_cons$X3)
food_cons$yi=food_cons$Y-mean(food_cons$Y)
```

Thus,the data in deviation form is shown in Table 2.3.

The model is:

```
#####
# Total
fit=lm(Y~X2+X3,data=food_cons)
#summary(fit)$coef
```

The estimated model is Table 2.4:

$$\hat{Y} = 86.0832 - 0.216X_2 + 0.3781X_3$$

To estimate variances of β 's, first determine $\hat{\sigma}^2$

```
## Var(b3_hat)
rsdl<-resid(fit)
(sigma_hat<-sum(rsdl^2)/(nrow(food_cons)-3)) #k=3

## [1] 0.7139
```

The estimated errors (residuals) are:

$$\hat{\mathcal{E}}_i = y_i - \hat{Y} = y_i - 86.0832 - 0.216X_2 - 0.3781X_3$$

The error sum of squares (ESS) = $\sum \hat{\mathcal{E}}_i^2 = 8.5673$

An estimator of the error variance σ^2 is:

$$\hat{\sigma}^2 = \frac{\sum \hat{\mathcal{E}}_i^2}{n-3} = \frac{8.5673}{12} = 0.7139$$

Estimation of standard errors of estimated coefficients

The standard errors of estimated regression coefficients for $\hat{\beta}_2$ is estimated as using equation (2.28) and (2.29) as follow:

First regress X2 on the others X's, and obtain the correlation

```
#####
fit_x2<-lm(X2~X3,data = food_cons)
#The R_s for regressing X2 on the other X's is given by:
#R2=1-ESS/TSS
#k=2
# Residual for regressing X3 on other X's(, that is in this case X2)
v2<-resid(fit_x2)
#R2=1-(ESS/TSS)
(R_sq_2<-1-(sum(v2^2)/sum(food_cons$x2i^2)))
```

```
## [1] 0.2557
```

```
#var_b2=sigma_hat/()  
var_b2=sigma_hat/sum(v2^2)  
(s.e_b2=sqrt(var_b2))
```

```
## [1] 0.05197
```

```
#OR sum(v2^2)=ESS(1-R_sq)  
var_b2=sigma_hat/(sum(food_cons$x2i^2)*(1-R_sq_2))  
(s.e_b2=sqrt(var_b2))
```

```
## [1] 0.05197
```

$$s.e.(\hat{\beta}_2) = 0.052$$

Similarly, the standard errors of estimated regression coefficients for $\hat{\beta}_3$ is estimated as:

First regress X_3 on the others X 's, and obtain the correlation

```
fit_x3<-lm(X3~X2,data = food_cons)  
#The R_s for regressing X3 on the other X's is given by:  
#R2=1-ESS/TSS  
#k=2  
# Residual for regressing X3 on other X's(, that is in this case X2)  
v3<-resid(fit_x3)  
#R2=1-(ESS/TSS)  
R_sq_3<-1-(sum(v3^2)/sum(food_cons$x3i^2))  
#var_b3=sigma_hat/()  
var_b3=sigma_hat/sum(v3^2)  
(s.e_b3=sqrt(var_b3))
```

```
## [1] 0.03383
```



```
#OR sum(v3^2)=ESS(1-R_sq)
var_b3=sigma_hat/(sum(food_cons$x3i^2)*(1-R_sq_3))
(s.e_b3=sqrt(var_b3))

## [1] 0.03383
```

$$s.e.(\hat{\beta}_3) = 0.0338$$

The coefficient of determination is given by:

```
summary(fit)$r.squared

## [1] 0.9143
```

$$R^2 = 0.914$$

1. $R^2 = 0.914$ indicates that 91.4% of the variation (change) in food consumption is attributed to the effect of food price and consumer income.
2. $1 - R^2 = 0.086$. This indicates that 8.6% of the variation in food consumption is due to factors (variables) not included in our specification.

Tests of model adequacy

A test of model adequacy is accomplished by testing the null hypothesis:

$$H_0 : \beta_2 = \beta_3 = 0$$

$$H_A : H_0 \text{ is not true}$$

The test statistic for this test is given by:

$$F_{cat} = \frac{RSS/(k-1)}{ESS/(n-k)} = \frac{91.362/(3-1)}{8.567271/(15-3)} = 63.98448$$

We compare this F-ratio with $F_\alpha(K-1, n-K) = F_\alpha(2, 12)$ for some significance level α .

- For $\alpha=0.01$, $F_{\alpha}(K - 1, n - K) = F_{0.01}(2, 12)=6.93$

- For $\alpha= 0.05$, $F_{\alpha}(K - 1, n - K) = F_{0.05}(2, 12)=3.89$

Since the test statistic is greater than both tabulated values, the above ratio is significant at the conventional levels of significance (1% and 5%). Thus, we reject the null hypothesis and conclude that the model is adequate, that is, variation (change) in per capita food consumption is **significantly attributed to** the effect of food price and/or per capita disposable income.

Tests of significance of regression coefficients

a) Does food price significantly affect per capita food consumption?

The hypothesis to be tested is:

$$H_0 : \beta_2 = 0$$

$$H_A = \beta_2 \neq 0$$

The test statistic is calculated as:

$$t_2 = \frac{\hat{\beta}_2}{s.e.(\hat{\beta}_2)} = \frac{-0.21596}{0.05197} = -4.155$$

For significance level $\alpha = 0.01$ and degrees of freedom $(n - 3) = (15 - 3) = 12$, the value from the student's t-distribution is:

$$t_{\alpha/2}(n - 3) = t_{0.005}(12) = 3.055$$

Decision: Since $|t_2| = 4.155 > 3.055$, we reject the null hypothesis and conclude that food price significantly affects per capita food consumption at the 1% level of significance.

b) Does disposable income significantly affect per capita food consumption? The hypothesis to be tested is:

$$H_0 : \beta_3 = 0$$

$$H_A = \beta_3 \neq 0$$

The test statistic is calculated as:

$$t_3 = \frac{\hat{\beta}_3}{s.e.(\hat{\beta}_3)} = \frac{0.378127}{0.033826} = 11.179$$

The 1% critical value from the student's t-distribution is again 3.055.

Decision: Since $|t_3| = 11.179 > 3.055$, we reject the null hypothesis and conclude that disposable income significantly affects per capita food consumption at the 1% level of significance.

Generally we have the following:

- Food price significantly and **negatively** affects per capita food consumption, while disposable income significantly and **positively** affects per capita food consumption.
- The estimated coefficient of food price is -0.21596. Holding disposable income constant, a one dollar increase in food price results in a 0.216 dollar decrease in per capita food consumption.
- The estimated coefficient of food price is 0.378127. Holding food price constant, a one dollar increase in disposable income results in a 0.378 dollar increase in per capita food consumption.

As can be seen from Table 4, the p-values for price and income are both less than 0.01. Thus, we can conclude that both variables significantly affect consumption at the 1% level of significance.

From the signs of the estimated regression coefficients we can see that the direction of influence is opposite: price affects consumption negatively while income affects consumption positively. The constant term (intercept) is also significant.

Note: In general, if the p-value > 0.05 , then we doubt the importance of the variable!

2.2.4 Matrix form of the multiple linear regression model

Consider the model:

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k X_{ki} + \mathcal{E}_i$$

Since we have n observations, we can write the model for each observed value as:

$$Y_1 = \beta_1 + \beta_2 X_{21} + \beta_3 X_{31} + \dots + \beta_k X_{k1} + \mathcal{E}_1$$

$$Y_2 = \beta_1 + \beta_2 X_{22} + \beta_3 X_{32} + \dots + \beta_k X_{k2} + \mathcal{E}_2$$

.

.

.

$$Y_n = \beta_1 + \beta_2 X_{2n} + \beta_3 X_{3n} + \dots + \beta_k X_{kn} + \mathcal{E}_n$$

The matrix form of the above model is:

$$Y = X\beta + \mathcal{E}$$

where

$$Y_{(n \times 1)} = \begin{bmatrix} Y_1 \\ Y_2 \\ \cdot \\ \cdot \\ \cdot \\ Y_n \end{bmatrix}, \quad X_{(n \times k)} = \begin{bmatrix} 1 & X_{21} & \cdot & \cdot & \cdot & X_{k1} \\ 1 & X_{22} & \cdot & \cdot & \cdot & X_{k2} \\ \cdot & \cdot & & & & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & & & & \cdot \\ 1 & X_{2n} & \cdot & \cdot & \cdot & X_{kn} \end{bmatrix}, \quad \beta_{(k \times 1)} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \cdot \\ \cdot \\ \cdot \\ \beta_n \end{bmatrix}, \quad \text{and} \quad \mathcal{E}_{(n \times 1)} = \begin{bmatrix} \mathcal{E}_1 \\ \mathcal{E}_2 \\ \cdot \\ \cdot \\ \cdot \\ \mathcal{E}_n \end{bmatrix}$$

Note:

- a) One of the assumptions of the classical linear regression model is that there is no multicollinearity, meaning that there is no linear dependence between the regressor variables

X_2, X_3, \dots, X_k . This is the same as saying that the matrix X has **full columnrank**. Since X has K columns, we write this as:

$$\text{Rank} (X) = K$$

- b) If a matrix X is of full rank K , then $X'X$ is also of full rank, i.e., $\text{Rank} (X'X) = K$. In such cases the inverse of $X'X$ exists. Otherwise (that is, if $\text{Rank} (X'X) < K$), then $X'X$ is said to be **singular**(and its inverse does not exist).

The **mean** of the error vector \mathcal{E} is:

$$E(\mathcal{E}) = \begin{bmatrix} E(\mathcal{E}_1) \\ E(\mathcal{E}_2) \\ \cdot \\ \cdot \\ \cdot \\ E(\mathcal{E}_n) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \cdot \\ \cdot \\ \cdot \\ 0 \end{bmatrix} = 0$$

Definition :

If Y is an $(n \times 1)$ random vector, then the variance-covariance matrix of Y is given by:

$$\sum_Y = E \{ [Y - E(Y)][Y - E(Y)]' \}$$

The variance-covariance matrix of $\hat{\beta}$ is:

$$\begin{aligned}
var(\hat{\beta}) &= E \left\{ [\hat{\beta} - E(\hat{\beta})][\hat{\beta} - E(\hat{\beta})'] \right\} \\
&= E \left\{ [\hat{\beta} - \beta][\hat{\beta} - \beta]' \right\}, (since E(\hat{\beta}) = \beta) \\
&= E \left\{ [(X'X)^{-1}X'\mathcal{E}][(X'X)^{-1}X'\mathcal{E}]' \right\} (from(*)) \\
&= E \left\{ [(X'X)^{-1}X'\mathcal{E}][\mathcal{E}'X(X'X)^{-1}] \right\} \\
&= (X'X)^{-1}X' \underbrace{E(\mathcal{E}\mathcal{E}')} X(X'X)^{-1}, \dots (where, E(\mathcal{E}\mathcal{E}') = \sigma^2 I_k) \\
&= \sigma^2 (X'X)^{-1} \underbrace{X'X(X'X)^{-1}} \dots (where X'X(X'X)^{-1} = I_k) \\
&= \sigma^2 (X'X)^{-1}
\end{aligned}$$

Estimation of σ^2

An unbiased estimator of σ^2 is given by:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{\mathcal{E}}_i^2}{n - k} = \frac{\hat{\mathcal{E}}'\hat{\mathcal{E}}}{n - k}$$

where $\hat{\mathcal{E}} = Y - X\hat{\beta}$. Thus, an estimator of the variance-covariance matrix of $\hat{\beta}$ is

$$\hat{\sum}_{\hat{\beta}} = \hat{\sigma}^2 (X'X)^{-1} = \frac{\hat{\mathcal{E}}'\hat{\mathcal{E}}}{n - k} (X'X)^{-1}$$

For instance, for case of k=2 is given by:

$$\text{var}(\hat{\beta}) = \hat{\sigma}^2(X'X)^{-1}$$

$$\begin{aligned}
&= \hat{\sigma}^2 \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ x_{11} & x_{12} & x_{13} & \dots & x_{1n} \end{bmatrix} \begin{bmatrix} 1 & x_{11} \\ 1 & x_{12} \\ 1 & x_{13} \\ \cdot & \cdot \\ \cdot & \cdot \\ 1 & x_{1n} \end{bmatrix} = \hat{\sigma}^2 \begin{bmatrix} n & \sum_{i=1}^n x_{1i} \\ \sum_{i=1}^n x_{1i} & \sum_{i=1}^n x_{1i}^2 \end{bmatrix}^{-1} \\
&= \frac{\hat{\sigma}^2}{n \times \sum_{i=1}^n x_{1i}^2 - (\sum_{i=1}^n x_{1i})^2} \begin{bmatrix} \sum_{i=1}^n x_{1i}^2 & -\sum_{i=1}^n x_{1i} \\ -\sum_{i=1}^n x_{1i} & n \end{bmatrix} \\
\text{cov}(\hat{\alpha}, \hat{\beta}) &= \frac{-\hat{\sigma}^2 \sum_{i=1}^n x_{1i}}{n \times \sum_{i=1}^n x_{1i}^2 - (\sum_{i=1}^n x_{1i})^2} \\
\text{Use : } n \sum_{i=1}^n x_{1i}^2 - (\sum_{i=1}^n x_{1i})^2 &= n \sum x^2
\end{aligned}$$

Note that $\hat{\Sigma}_{\hat{\beta}}$ is a $(K \times K)$ matrix. The main diagonal elements of this matrix are the estimated variances of $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$. The square roots of the main diagonal entries of $\hat{\Sigma}_{\hat{\beta}}$ are the standard errors: $s.e.(\hat{\beta}_1), s.e.(\hat{\beta}_2), \dots, s.e.(\hat{\beta}_k)$.

To test for the significance of the regression coefficients:

$$H_0 : \beta_J = 0$$

$$H_A : \beta_J \neq 0$$

the test statistic is:

$$t_j = \frac{\hat{\beta}_j}{s.e.(\hat{\beta}_j)}, j = 1, 2, \dots, k$$

Decisionrule :

If $|t_j| > t_{\alpha/2}(n - K)$, we reject H_0 and conclude that β_j is significant, that is, the regressor variable $X_j, j=1, 2, \dots, K$, significantly affects the dependent variable Y .

2.3 Heteroscedasticity

In this topic we will see

- the consequences of violation of constant variance assumption.
- What are the consequences for the properties of least squares estimators?
- Is there a better estimation technique?
- How do we detect the existence of Heteroscedasticity?

2.3.1 Introduction

Consider the relationship between average or mean household expenditure on food $E(y)$ and household income x described by the linear function:

$$E(y) = \alpha + \beta x$$

The unknown parameters α and β convey information about this expenditure function. The **response parameter** β describes how mean household food expenditure changes when household income increases by one unit. The intercept parameter α measures expenditure on food for a zero income level. Knowledge of these parameters aids planning by institutions such as government agencies or food retail chains.

To estimate α and β , use data **food.csv** and this data consists a sample of 40 households. Let \mathcal{E}_i be the difference between expenditure on food by the i^{th} household y_i and mean expenditure on food for all households with income x_i . That is,

$$\mathcal{E}_i = y_i - E(y_i) = y_i - \alpha - \beta x_i$$

$$\implies y_i = \alpha + \beta x_i + \mathcal{E}_i$$

Model used to describe expenditure on food for the i^{th} household is written as

This can be viewed $E(y_i) = \alpha + \beta x_i$ as that part of food expenditure explained by income x_i and \mathcal{E}_i as that part of food expenditure explained by other factors.

Now, we need to check from the data: Whether the mean function $E(y_i) = \alpha + \beta x_i$ is **equal or better** at explaining expenditure on food for low-income households than it is for high-income households or not?.

As we know that low-income households do not have the option of extravagant food tastes.

Comparatively, they have few choices and are almost forced to spend a particular portion of their income on food. High-income households on the other hand could have simple food tastes or extravagant food tastes. Thus, income is relatively **less important** as an explanatory variable for food expenditure of high-income households.

It is harder to guess their food expenditure. This can also be described as that the probability of getting large positive or negative values for \mathcal{E} is higher for high incomes than it is for low incomes. Factors other than income can have a larger impact on food expenditure when household income is high. So, how can we model this phenomenon? A random variable, in this case \mathcal{E} , has a higher probability of taking on large values if its variance is high. Thus, we can capture the effect we are describing by having $Var(\mathcal{E})$ depends directly on income x . An equivalent statement is to say $Var(y)$ increases as x increases. Food expenditure y can deviate further from its mean $E(y) = \alpha + \beta x$ when x is large. In such a case, when the variances for all observations are not the same, we say that heteroskedasticity exists (see Figure 2.7). Alternatively, we say the random variable y and the random error \mathcal{E} are heteroskedastic. Conversely, if all observations come from probability density functions with the same variance, we say that homoskedasticity exists, and y and \mathcal{E} are homoskedastic.

From Figure 2.7, at $x = x_1$, the probability density function $f(y_1|x_1)$ is such that y_1 will be close to $E(y_1)$ with high probability. When we move to x_2 , the probability density function $f(y_2|x_2)$ is more spread out; we are less certain about where y_2 might fall, and larger values are possible.

One of the assumptions of CLRM is that random error terms are uncorrelated with mean zero and

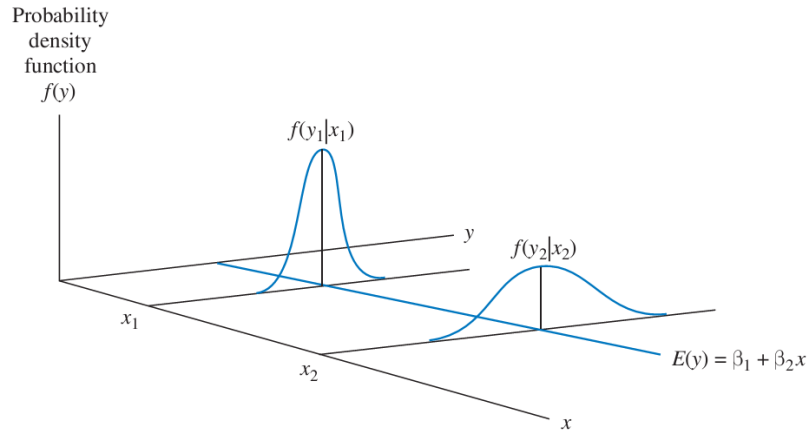


Figure 2.7: Heteroskedastic errors.

constant variances σ^2 :

$$Var(\mathcal{E}_i) = E(\mathcal{E}_i^2) = \sigma^2, (i = 1, 2, 3, \dots, n)$$

This assumption tells us that the variance remains constant for all observations. But there are many situations as we saw in the above in which this assumption may not hold. For example, the variance of the error term may increase or decrease with the dependent variable or one of the independent variables. Under such circumstances, we have the case of **heteroscedasticity**. Generally, under heteroscedasticity we have:

$$E(\mathcal{E}_i^2) = k_i \sigma^2, k_i \text{ are not all equal or}$$

$$Var(y_i) = Var(\mathcal{E}_i) = g(x_i)$$

The least squares estimated equation from **food.csv** data is shown in Figure 2.8

```
fit<-lm(food_exp~income,data=food)

plot(food$income, food$food_exp, pch=19, xlab = "x=weekly income in $100", ylab=
```

```
points(food$income,fitted(fit),type="l",lwd=2)
#text(10,500,expression(hat(y)==83.42 + 10.21x))
text(15,500,expression(paste(hat(y)==83.42 + 10.21,"x")))
```

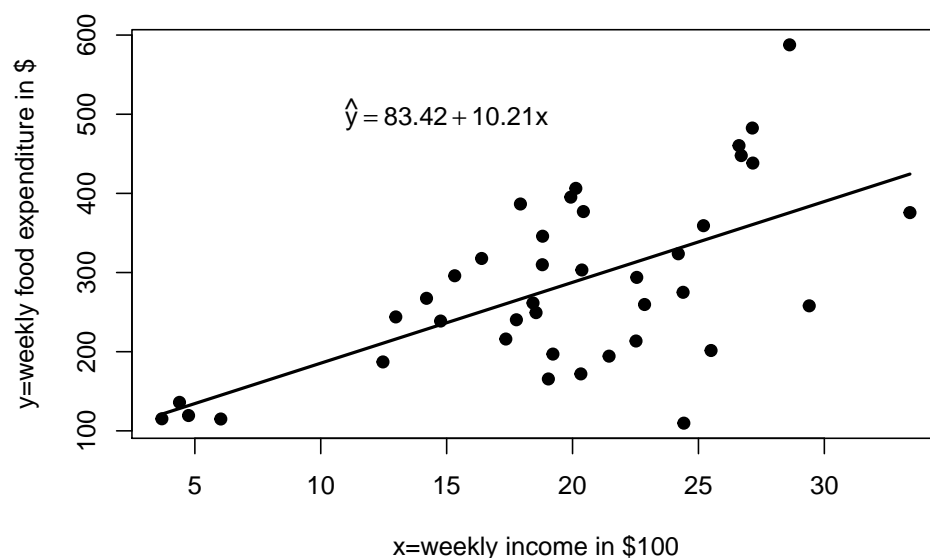


Figure 2.8: Least squares estimated food expenditure function and observed data points

Heteroscedasticity is more plausible at cross-sectional data than time series data. Here the assumption of homoscedasticity is not very plausible since we expect less variation in consumption for low income families than for high income families. At low levels of income, the average level of consumption is low and the variation around this level is restricted: consumption can not fall too far below the average level because this might mean starvation, and it can not rise too far above the average level because the asset does not allow it. These constraints are likely to be less binding at higher income levels.

2.3.2 Consequences of heteroscedasticity

Since the existence of heteroskedasticity means that the least squares assumption $var(\mathcal{E}_i) = \sigma^2$ is violated, we need to ask what consequences this violation has for our least squares estimator, and what we can do about it.

There are two implications:

- 1) The least squares estimator is still a linear and unbiased estimator, but it is no longer best. There is another estimator with a smaller variance.
- 2) The standard errors usually computed for the least squares estimator are incorrect (estimated variances of the OLS estimators are biased). Confidence intervals and hypothesis tests that use these standard errors may be misleading.

Consider the simple linear regression model (in deviation form):

$$y_i = \beta x_i + \mathcal{E}_i, (i = 1, 2, 3, \dots, n)$$

where \mathcal{E}_i satisfies all assumptions of the CLRM except that the error terms are heteroscedastic, that is,

$$E(\mathcal{E}_i^2) = k_i \sigma^2 \text{ and } k_i \text{ are not all equal.}$$

The OLS estimator of β is:

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

Then we have:

$$\begin{aligned} \hat{\beta} &= \frac{\sum x_i y_i}{\sum x_i^2} = \beta \frac{\sum x_i^2}{\sum x_i^2} + \frac{\sum x_i \mathcal{E}_i}{\sum x_i^2} = \beta + \frac{\sum x_i \mathcal{E}_i}{\sum x_i^2} \\ \Rightarrow E(\hat{\beta}) &= \beta + \frac{\sum x_i E(\mathcal{E}_i)}{\sum x_i^2} = \beta, \text{ where } E(\mathcal{E}_i) = 0 \end{aligned}$$

Thus, $\hat{\beta}$ is an unbiased estimator of β even in the presence of heteroscedasticity.

Recall that the variance of the OLS estimator $\hat{\beta}$ when there is no heteroscedasticity (or under homoscedasticity) is given by:

$$Var(\hat{\beta}) = \frac{\sigma^2}{\sum x_i^2} \quad (2.32)$$

Under heteroscedasticity we have:

$$\begin{aligned}
Var(\hat{\beta})_{HET} &= E(\hat{\beta} - E(\hat{\beta}))^2 = E(\hat{\beta} - \beta)^2 \\
&= (-1)^2 E(\beta - \hat{\beta})^2 \quad \dots\dots\dots(Unbiasedness) \\
&= E\left(\frac{\sum x_i \mathcal{E}_i}{\sum x_i^2}\right)^2 \\
&= E\left(\frac{\sum x_i^2 \mathcal{E}_i^2}{(\sum x_i^2)^2} + \frac{\sum_{i \neq j} x_i x_j \mathcal{E}_i \mathcal{E}_j}{(\sum x_i^2)^2}\right) \\
&= \frac{\sum x_i^2 E(\mathcal{E}_i^2)}{(\sum x_i^2)^2} + \frac{\sum_{i \neq j} x_i x_j E(\mathcal{E}_i \mathcal{E}_j)}{(\sum x_i^2)^2}, \\
&\quad (where E(\mathcal{E}_i^2) = k_i \sigma^2 \text{ and } E(\mathcal{E}_i \mathcal{E}_j) = 0) \\
&= \frac{\sigma^2 \sum k_i x_i^2}{(\sum x_i^2)^2} = \frac{\sigma^2}{\sum x_i^2} \frac{\sum k_i x_i^2}{\sum x_i^2} \\
&= Var(\hat{\beta}) \left[\frac{\sum k_i x_i^2}{\sum x_i^2} \right] \tag{2.33}
\end{aligned}$$

$$\begin{aligned}
&\text{if you assume } var(\mathcal{E}_i) = \sigma_i^2 \\
&= \frac{\sum (X_i - \bar{x})^2 \sigma_i^2}{[\sum (X_i - \bar{x})^2]^2} \\
&= \frac{\sum (x_i)^2 \sigma_i^2}{[\sum (x_i)^2]^2} \tag{2.34}
\end{aligned}$$

Consequently, if we proceed to use the least squares estimator and its usual standard errors when $Var(\mathcal{E}_i) = k_i \sigma^2 = \sigma_i^2$, we are committing **two errors**.

- 1) we are using an estimate of equation (2.32) to compute the standard error of $\hat{\beta}$ when we should be using an estimate of (2.33).
- 2) it is using s^2 to estimate a common σ^2 when in fact the σ_i^2 's are different.

From (2.32) and (2.33), it can be seen that the two variances, $Var(\hat{\mathcal{E}})$ and $Var(\hat{\mathcal{E}})_{HET}$ will be equal only if $k_i = 1 \forall_i$ that is, only if the errors are homoscedastic.

- 1) If $\frac{\sum k_i X_i^2}{\sum X_i^2} < 1$, then OLS will overestimate the variance of $\hat{\beta}$.
- 2) If $\frac{\sum k_i X_i^2}{\sum X_i^2} > 1$, then OLS will underestimate the variance of $\hat{\beta}$.

Thus, under heteroscedasticity, the OLS estimators of the regression coefficients are not BLUE and efficient(with minimum variance). We have to find an alternative estimator that has the minimum

variance property.

2.3.3 Detecting Heteroskedasticity

How do I know if heteroskedasticity is likely to be a problem for my model and my set of data? Is there away of detecting heteroskedasticity so that I know whether to investigate other estimation techniques? We consider three ways of investigating these questions.

The tests consider a test for heteroskedasticity based on a **variance function**. As we know that the mean function $E(y_i)$ for a general multiple linear regression is given by:

$$E(y_i) = \beta_1 + \beta_2 x_{2i} + \dots + \beta_k x_{ki}$$

A general form for the variance function also given by:

$$var(y_i) = \sigma_i^2 = E(\mathcal{E}_i^2) = h(\alpha_1 + \alpha_2 z_{2i} + \dots + \alpha_s z_{si})$$

1. Residual Plots

One way of investigating the existence of heteroskedasticity is to estimate your model using least squares and to plot the least squares residuals. If the errors are heteroskedastic, they may tend to exhibit greater variation in some systematic way.

In a regression with more than one explanatory variable we can plot the least squares residuals against each explanatory variable, or against \hat{y}_i , to see if they vary in a systematic way.

Example

Consider the data on weekly food expenditure (Y) and weekly income in \$100(X) for 20 households(both in thousands of Dollars)(use “**consmpction.csv**”).

```
consmpction<-read.csv("consmpction.csv",header=T,sep=",")
head(consmpction,1)
```

```
fit<-lm(expen~income,data=consmpction)
#plot(fit$fitted.values,fit$residuals,xlab="x=fitted",ylab="Residuals")
```

```
#abline(h=0,lwd=2,col=4)
# or use
plot(consmption$income,fit$residuals,xlab = "x=income",ylab = "Residuals")
abline(h=0,lwd=2,col="red")
```

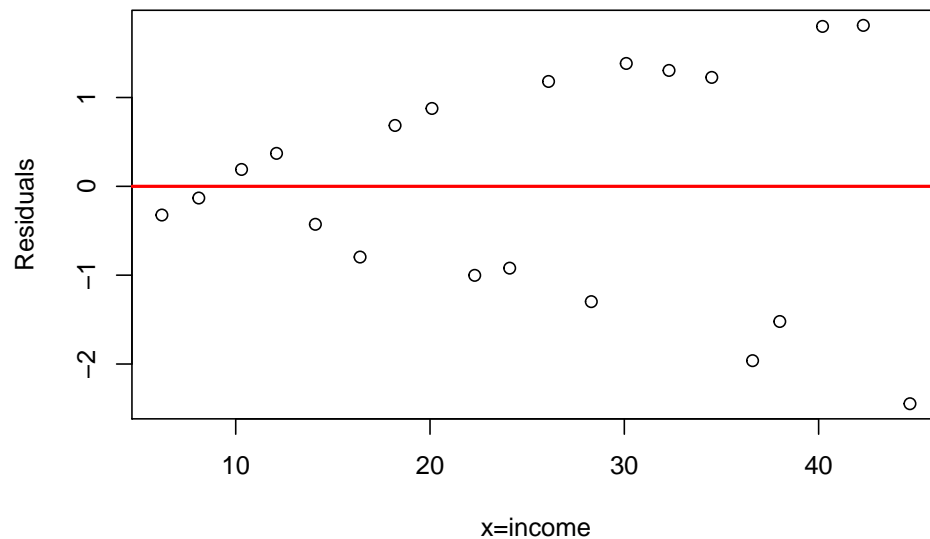


Figure 2.9: Residual Plot

From the graph, we suspect that the variance increases as incomes increases.

2. White's test

It is a general test for homoskedasticity where **nothing is known** about the form of this heteroskedasticity. This test is based on the difference between the variance of the OLS estimates under homoskedasticity and that under heteroskedasticity.

A test for heteroskedasticity without precise knowledge of the relevant variables and instead by defining the z 's as equal to the x 's, the squares of the x 's, and possibly their cross-products. This test involves applying OLS to:

$$\hat{\mathcal{E}}_i^2 = \gamma_0 + \gamma_1 Z_{1i} + \gamma_2 Z_{2i} + \dots + \gamma_p Z_{pi} + U_i$$

and calculating the coefficient of determination R_w^2 , where $\hat{\mathcal{E}}_i$ are OLS residuals from the original model. The null hypothesis is:

$$H_0 : \gamma_1 = \gamma_2 = \dots = \gamma_p = 0$$

The test statistic is:

$$\chi_{cal}^2 = nR_w^2$$

Decisionrule : Reject H_0 (the hypothesis of homoscedasticity) if the above test statistic exceeds the value from the Chi-square distribution with p degrees of freedom for a given level of significance α .

If our model has only one independent variable X_i , then $p = 2$: $Z_{1i} = X_i$ and $Z_{2i} = X_i^2$. If we have two independent variables X_{1i} and X_{2i} , then $p = 5$: $Z_{1i} = X_{1i}$, $Z_{2i} = X_{2i}$, $Z_{3i} = X_{1i}X_{2i}$, $Z_{4i} = X_{1i}^2$ and $Z_{5i} = X_{2i}^2$. And so on.

Example

For the consumption data, test heteroskedasticity? (use “**consmpcion.csv**”).

solution

To test for heteroskedasticity in the food expenditure where the variance is potentially a function of income, we test $H_0 : \delta_1 = \delta_2 = 0$ Vs $H_1 : \delta_1 \neq 0$ or $\delta_2 \neq 0$ in the variance function

$$\sigma_i^2 = h(\alpha_1 + \alpha_2 x_i).$$

This involves applying OLS to: $\hat{\mathcal{E}}_i^2 = \delta_0 + \delta_1 X_i + \delta_2 X_i^2 + U_i$ and computing the coefficient of determination R_w^2 .

```
consmpcion<-read.csv("consmpcion.csv",header=T,sep=",")
head(consmpcion,1)
```

```
fit<-lm(expen~income,data=consmpcion)
consmpcion$ei_hat<-resid(fit)
fit_w<-(lm(I(ei_hat^2)~income+I(income^2),data=consmpcion))
(Rw_square<-summary(fit_w)$r.squared)
```



```
## [1] 0.8781
```

```
x_2<-nrow(consmption)*Rw_square
```

This yields $R_w^2 = 0.8781$. The White test statistic is:

$$\chi_{cal}^2 = nR_w^2 = 20(0.8781) = 17.5623$$

We compare this value with $\chi_\alpha^2(p)$ for a given level of significance α . For $\alpha = 0.05$,

```
(chi_tab<-qchisq(0.05,2,lower.tail=F))
```

```
## [1] 5.991
```

```
# Or using a p-value<0.05
```

```
pchisq(q=x_2,df=1,lower.tail=F)
```

```
## [1] 2.78e-05
```

$\chi_{0.05}^2(2) = 5.9915$.

Conclusion: we conclude that heteroskedasticity exists with the variance dependent on income.

3. Goldfeld-Quandt test

This test for heteroskedasticity is designed for two groups of data with possibly different variances. To introduce this case, consider a wage equation where earnings per hour (WAGE) depends on years of education (EDUC), years of experience (EXPER) and a dummy variable METRO that is equal to one for workers who live in a metropolitan area and zero for workers who live outside a metropolitan area. Using data in the file **cps2.csv**(metro= 1 if lives in metropolitan area,female= 1 if female,).

```
cps2<-read.csv("cps2.csv",header = T,sep = ",")  
head(cps2,2)
```

```
fit<-lm(wage~educ+exper+metro,data=cps2)
```

The least squares estimated equation for this model is

$$WAGE = -9.914 + 1.234EDUC + 0.1332EXPER + 1.5241METRO$$

The results suggest that education and experience have a positive effect on the level of wages and that given a particular level of education and experience, the average metropolitan wage is \$1.50 per hour higher than the average wage in a rural area.

The question we now ask is: 'How does the variance of wages in a metropolitan area compare with the variance of wages in a rural area? Are the variances likely to be the same, or different? One might suspect that the greater range of different types of jobs in a metropolitan area might lead to city wages' having a higher variance. If the variance of metropolitan wages differs from the variance of rural wages, then we have heteroskedasticity. The variance is not constant for all observations. The Goldfeld-Quandt test is designed to test for this form of heteroskedasticity, where **the sample can be partitioned into two groups**—metropolitan and rural in this case—and we suspect the variance could be different in the two groups.

Suppose we have a model with one explanatory variable X_1 and let Y be the dependent variable. Assuming that the variance is related with X_1 and is continuous variable(). The steps involved in this test are the following:

- a) Arrange the observations (both Y and X_1) in increasing order of X_1 .
- b) Divide the observations into three parts: n_1 observations in the first part, p observations in the middle part, and n_2 observations in the second part ($n_1 + n_2 + p = n$). Usually p is taken to be **one-sixth** of n .
- c) Run a regression on the first n_1 observations, obtain the residuals $\hat{\mathcal{E}}_{1i}$, and calculate the residual variance $s_1^2 = \sum_{i=1}^{n_1} \frac{\hat{\mathcal{E}}_{1i}^2}{(n_1-2)}$. Similarly run a regression on the second n_2 observations, obtain the residuals $\hat{\mathcal{E}}_{2i}$, and calculate the variance $s_2^2 = \sum_{i=1}^{n_2} \frac{\hat{\mathcal{E}}_{2i}^2}{(n_2-2)}$.

Note:

The variances of the last several disturbances in the first part are likely to be similar to those of the first several disturbances in the second part. To increase the power of the test, it is recommended that the two parts be some distance apart. Thus, we drop the middle p residuals all together.

If we need to test variance for a categorical variable like (Urban/Rural), directly apply **step c**

separately.

d) Calculate the test statistic: $F_{cal} = \frac{S_2^2}{S_1^2}$

e) *Decisionrule*: Reject the null hypothesis $H_0 : \sigma_1^2 = \sigma_2^2$ (and conclude that the errors are heteroscedastic) if:

$$F_{cal} > F_{\alpha}(n_1 - 2, n_2 - 2)$$

where $F_{\alpha}(n_1 - 2, n_2 - 2)$ is the critical value from the F-distribution with $n_1 - 2$ and $n_2 - 2$ degrees of freedom in the numerator and denominator, respectively, for a given significance level α .

Example

Apply the Goldfeld-Quandt test for non-constant error variance for metropolitan area and non-metropolitan area(rural)? use (**cps2.csv**).

solution

The test is based on a comparison of the error variances estimated from each group. Using the subscript M to denote metropolitan observations and the subscript R to denote rural observations, we can write separate equations for the two groups as

$$WAGE_{Mi} = \beta_{M1} + \beta_2 EDUC_{Mi} + \beta_3 EXPER_{Mi} + \mathcal{E}_{Mi} \quad (2.35)$$

$$WAGE_{Ri} = \beta_{R1} + \beta_2 EDUC_{Ri} + \beta_3 EXPER_{Ri} + \mathcal{E}_{Ri} \quad (2.36)$$

Implicit in the above specification is the assumption that the coefficients for EDUC and EXPER (β_2 and β_3) are the same in both metropolitan and rural areas, but the intercepts differ.

Now we need to test $H_0 : \sigma_M^2 = \sigma_R^2$ Vs $H_1 : \sigma_M^2 \neq \sigma_R^2$

Apply **step c**: to both equation

```
fit_M<-lm(wage~educ+exper, data=cps2[cps2$metro==1,])
fit_R<-lm(wage~educ+exper, data=cps2[cps2$metro==0,])
(s_M_sq<-sum(resid(fit_M)^2)/(nrow(cps2[cps2$metro==1,])-3))
```

```
## [1] 31.82
```

```
(s_R_sq<-sum(resid(fit_R)^2)/(nrow(cps2[cps2$metro==0,])-3))
```

```
## [1] 15.24
```

```
(f_cal<-s_M_sq/s_R_sq)
```

```
## [1] 2.088
```

Calculate

$$F_{cal} = \frac{s_M^2}{s_R^2} = 31.8237/15.243 = 2.0878$$

The lower and upper critical values for a 5% significance level are

$F_{Lc} = F_{(0.025,805,189)}$ and $F_{Uc} = F_{(0.975,805,189)}$ are given below. We reject

H_0 if $F_{cal} < F_{Lc}$ or $F_{cal} > F_{Uc}$

```
(F_Lc<-qf(0.025,nrow(cps2[cps2$metro==1,])-3,nrow(cps2[cps2$metro==0,])-3))
```

```
## [1] 0.8052
```

```
(F_Uc<-qf(0.975,nrow(cps2[cps2$metro==1,])-3,nrow(cps2[cps2$metro==0,])-3))
```

```
## [1] 1.262
```

Conclusion: Since $2.0878 > 1.2617$, we reject H_0 and conclude that the wage variances for the rural and metropolitan regions are not equal.

4. Breusch-Pagan test

This involves applying OLS to:

$$\frac{\hat{\mathcal{E}}_i^2}{\hat{\sigma}^2} = \gamma_0 + \gamma_1 X_{1i} + \gamma_2 X_{2i} + \dots + \gamma_k X_{ki} + U_i$$

where $\hat{\sigma}^2 = \sum \hat{\mathcal{E}}_i^2/n$, the ML estimator of σ^2 under homoskedasticity.

and calculating the regression sum of squares (RSS). The test statistic is:

$$\chi_{cal}^2 = \frac{RSS}{2}$$

Decisionrule : Reject the null hypothesis of homoscedasticity: $H_0 : \gamma_1 = \gamma_2 = \dots = \gamma_k = 0$ if:

$$\chi_{cal}^2 > \chi_{\alpha}^2(K)$$

where $\chi_{\alpha}^2(K)$ is the critical value from the Chi-square distribution with K degrees of freedom for a given value of α .

Example

use data of **consumption.csv**, and test for presence of heteroscedasticity using Breusch - Pagan test?

solution

This involves applying OLS to:

$$\frac{\hat{\mathcal{E}}_i^2}{\hat{\sigma}^2} = \gamma_0 + \gamma_1 X_i + U_i$$

```
#You have to know location of your dependent variable,
# when you use this function for other
bp_test<-function(x=data.frame(),alpha=null)
{
  n=nrow(x)
  fit<-lm(expen~income,data=x)
  sigma_hat_2=deviance(fit)/n
  res_bp<-resid(fit)^2/sigma_hat_2
  fit_bp<-lm(res_bp~income, data=x)
  fit_bp
  ## Regresssion SS
  rss=sum((fitted(fit_bp)-mean(res_bp))^2)
  # x_2 test statistic
  x_cal<-rss/2

  #Probability of getting greater than this calculated value,
  # that is, is this due to sampling error or not.
```

```

# p-value
p_value<-pchisq(q=x_cal, df=1,lower.tail=F)
#x_tab
x_tab<-qchisq(alpha, df=1,lower.tail=F)
# I recommend to use z next 11 commented codes instead of the last .
# if(p_value < alpha){
# cat("Decision : ", "\nReject Ho: constant variance, since p-value ",
#   # p_value, "<", alpha, "\nOr",
#   # "\nchi_tabulated=", x_tab, "<", x_cal, "=chi_calculated")
# }
# else{
# cat(" Do not reject Ho: constant variance, with a p-value", p_value)
# cat("Decision : ", "\nDo not reject Ho: constant variance, since p-value ",
#   #p_value, ">", alpha, "\nOr",
#   # "\nchi_tabulated=", x_tab, ">", x_cal, "=chi_calculated")
# }

if(p_value < alpha){
cat("Decision : ", "\nReject Ho: since p-value ", p_value, "<", alpha, "\nOr",
"\nchi_tabulated=", x_tab, "<", x_cal, "=chi_calculated")
}
else{
cat(" Do not reject Ho: since p-value", p_value)
cat("Decision : ", "\nDo not reject Ho:  since p-value ", p_value, ">", alpha, "
"\nX_tab=", x_tab, ">", x_cal, "=X_cal")
}
}

bp_test(x=consmpction, alpha = 0.05)

```

```
## Decision :
## Reject Ho:since p-value 0.006218 < 0.05
## Or
## chi_tabulated= 3.841 < 7.486 =chi_calculated
```

All of the tests indicated that the disturbances are heteroscedastic. Thus, the regression coefficients obtained by OLS are **not efficient**. In such cases, we have to apply another method such as weighted least squares (WLS) estimation.

Heteroskedasticity-Consistent Standard Errors

When our model suffers from heteroskedasticity, even the least squares estimator are unbiased but no longer best. The other is that the usual least squares standard errors are incorrect, which invalidates interval estimates and hypothesis tests. If we are prepared to accept the least squares estimator as a useful estimator, despite the fact it is not the minimum variance estimator, there is a way of correcting the standard errors so that our interval estimates and hypothesis tests are valid. Use a standard error of the estimator from equation (2.33) to test your hypothesis. These standard error are known by heteroskedasticity-consistent standard errors, or heteroskedasticity robust standard errors, or simply **robust standard errors**.

2.3.4 Correction for heteroscedasticity

Consider the model:

$$Y_i = \alpha + \beta X_i + \mathcal{E}_i \quad (2.37)$$

where $E(\mathcal{E}_i^2) = \sigma_i^2$ is known for $i = 1, 2, \dots, n$. We make the following transformation:

$$\frac{Y_i}{\sigma_i} = \frac{\alpha}{\sigma_i} + \beta \frac{X_i}{\sigma_i} + \frac{\mathcal{E}_i}{\sigma_i} \dots (\text{let, } \frac{Y_i}{\sigma_i} = Y_i^*, \frac{\alpha}{\sigma_i} = \alpha^*, \frac{X_i}{\sigma_i} = X_i^* \text{ and } \frac{\mathcal{E}_i}{\sigma_i} = \mathcal{E}_i^*).$$

The transformed model can be written as:

$$Y_i^* = \alpha^* + \beta X_i^* + \mathcal{E}_i^* \quad (2.38)$$

We then check if the disturbances of equation (2.38) satisfy OLS assumptions:

$$E(\mathcal{E}_i^*) = E\left(\frac{\mathcal{E}_i}{\sigma_i}\right) = \frac{E(\mathcal{E}_i)}{\sigma_i} = 0, \dots, (E(\mathcal{E}_i) = 0)$$

$$Var(\mathcal{E}_i^*) = E(\mathcal{E}_i^*)^2 = E\left(\frac{\mathcal{E}_i}{\sigma_i}\right)^2 = \frac{E(\mathcal{E}_i^2)}{\sigma_i^2} = E\left(\frac{\mathcal{E}_i}{\sigma_i}\right)^2 = \frac{\sigma_i^2}{\sigma_i^2} = 1$$

Thus, the variance of the transformed disturbances is constant. So we can apply OLS to equation (2.38) to get regression coefficient estimates that are BLUE (Gauss-Markov Theorem). This estimation method is known as **weighted least squares (WLS)** since each observation is weighted (multiplied) by $W_i = \frac{1}{\sigma_i}$.

Specification of the weights

The major difficulty with WLS is that σ_i^2 are rarely known. We can overcome this by making certain assumptions about σ_i^2 or by estimating σ_i^2 from the sample. The information about σ_i^2 is frequently in the form of an assumption that σ_i^2 is associated with some variable, say Z_i .

Illustration1 : In the case of micro-consumption function, the variance of the disturbances is often assumed to be positively associated with level of income. So the place of Z_i will be taken by the explanatory variable X_i income, that is,

$$\sigma_i^2 = \sigma^2 Z_i^2 = \sigma^2 X_i^2.$$

We then divide equation (2.37) throughout by X_i :

$$\begin{aligned} \frac{Y_i}{X_i} &= \alpha \left(\frac{1}{X_i} \right) + \beta \left(\frac{X_i}{X_i} \right) + \frac{\mathcal{E}_i}{X_i} \\ \Rightarrow \frac{Y_i}{X_i} &= \alpha \left(\frac{1}{X_i} \right) + \beta + u_i \end{aligned} \quad (2.39)$$

where $U_i = \frac{\mathcal{E}_i}{X_i}$. Now we have:

$$Var(U_i) = Var\left(\frac{X_i}{X_i}\right) = E\left(\frac{X_i}{X_i}\right)^2 = E\frac{\mathcal{E}_i^2}{X_i^2} = \frac{\sigma^2 X_i^2}{X_i^2} = \sigma^2$$

Hence, the variance of the disturbance term in equation (2.39) is constant, and we can apply OLS by regressing $\frac{Y_i}{x_i}$ on $\frac{1}{X_i}$. Note that the estimated constant term and slope from the transformed

model (2.39) correspond to the values of $\hat{\beta}$ and $\hat{\alpha}$, respectively.

Illustration 2: In the case of micro-consumption function, the variance of the disturbances may also be thought to be associated with changes in some ‘outside’ variable, say the size of the family= Z_i , that is,

$$\sigma_i^2 = \sigma^2 Z_i^2.$$

We then divide equation (2.36) throughout by Z_i :

$$\frac{Y_i}{Z_i} = \alpha \left(\frac{1}{Z_i} \right) + \beta \left(\frac{X_i}{Z_i} \right) + V_i$$

where $V_i = \frac{\varepsilon_i}{Z_i}$. It can easily be shown that $Var(V_i) = \sigma^2$. To estimate the regression coefficients, we should run a regression of $\frac{Y_i}{Z_i}$ on $\frac{1}{Z_i}$ and $\frac{X_i}{Z_i}$ without a constant term.

Example: apply WLS method on ‘consumption’ data?

All of the tests indicate that the disturbances are heteroscedastic. Thus, the regression coefficients obtained by OLS are not efficient. In such cases, we have to apply weighted least squares (WLS) estimation. The weights can be obtained from the

- **sample at hand** or
- from some **prior knowledge**.

In this example let’s estimate the weights σ_i from the sample (that is we are applying white’s test method).

First we apply OLS estimation and obtain the residuals $\hat{\varepsilon}_i$. We then order the residuals based on the absolute magnitude of the explanatory variable (income). Next we divide the residuals into three parts: the first and second parts consisting of seven residuals and the third part consisting of six residuals. The variance of each part is computed as:

$$\hat{\sigma}_i^2 = \frac{1}{n_i} \sum \hat{\varepsilon}_i^2$$

, where n_i is the number of residuals in the i^{th} part, $i = 1, 2, 3$.

The results are:

```

rm(list = ls())
consumption<-read.csv("consumption.csv",header=T,sep=",")
fit<-lm(expen~income,data=consumption)
sorted_residual<-resid(fit)[order(consumption$income)]
n1= 7; n2=7; n3=6
# variance of each part
sigma_hat1_sq<-sum(sorted_residual[1:7]^2)/n1
sigma_hat2_sq<-sum(sorted_residual[8:14]^2)/n2
sigma_hat3_sq<-sum(sorted_residual[15:20]^2)/n3
sigma_hat1_sq=round(sigma_hat1_sq,6)
sigma_hat2_sq=round(sigma_hat2_sq,6)
sigma_hat3_sq=round(sigma_hat3_sq,6)
s1=sqrt(sigma_hat1_sq); s2=sqrt(sigma_hat2_sq); s3=sqrt(sigma_hat3_sq)

```

$$\hat{\sigma}_1^2 = 0.2259 \implies \hat{\sigma}_1 = 0.4753$$

$$\hat{\sigma}_2^2 = 1.3306 \implies \hat{\sigma}_2 = 1.1535$$

$$\hat{\sigma}_3^2 = 3.363 \implies \hat{\sigma}_3 = 1.8339$$

The next step is to divide the values of the dependent variable, the independent variable and the constant term (a vector of 1's) in the i^{th} part by σ_i :

$$\frac{y_i}{\hat{\sigma}_i} = \alpha \left(\frac{1}{\hat{\sigma}_i} \right) + \beta \left(\frac{x_i}{\hat{\sigma}_i} \right) + U_i, \quad \text{where} \quad U_i = \frac{\mathcal{E}_i}{\hat{\sigma}_i}.$$

```

# Vector of weights
weight<-rep(c(s1,s2,s3),c(7,7,6))
xxx=consumption[, -1] # To remove Household variable

```

Table 2.5: transformed data divided by its respective sigma.

y_sigma_i	x_sigma_i	one_sigma_i
46.919	41.870	2.1040
67.959	65.645	2.1040
77.007	66.907	2.1040
25.459	25.459	2.1040
88.999	85.633	2.1040
13.045	12.834	2.1040
94.049	81.215	2.1040
22.626	22.106	0.8669
8.929	8.929	0.8669
34.850	33.636	0.8669

```
xxx$one<-1
weighted_data=xxx/weight
weighted_data<-round(weighted_data,4)
names(weighted_data)=c("y_sigma_i", "x_sigma_i", "one_sigma_i")
```

We then run an OLS regression of $\frac{1}{\sigma_i}$ on $\frac{X_i}{\sigma_i}$ and $\frac{Y_i}{\sigma_i}$ without a constant term. The results are:

```
# exclude the intercept (-1)
fit_wls<-lm(I(expen/weight)~I(1/weight)+I(income/weight)-1, consmption)
print(summary(fit_wls))
```

```
##
## Call:
## lm(formula = I(expen/weight) ~ I(1/weight) + I(income/weight) -
##      1, data = consmption)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.282 -0.619  0.036  0.967  4.621
##
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## I(1/weight)      0.9381      0.7585    1.24    0.23
## I(income/weight)  0.8881      0.0249   35.67   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.07 on 18 degrees of freedom
## Multiple R-squared:  0.997, Adjusted R-squared:  0.997
## F-statistic: 3.47e+03 on 2 and 18 DF, p-value: <2e-16
```

$$Y_i = \frac{0.9381}{(0.7585)} + \frac{0.8881X_i}{(0.0249)}$$

$$R^2 = 0.9974, \quad \hat{\sigma}^2 = 2.0668$$

The plot of the residuals of the transformed model against the explanatory variable (income) is shown below. It can be seen that the spread of the residuals has no increasing or decreasing pattern, i.e., there is no heteroscedasticity.

```
plot(consmption$income, fit_wls$residuals, xlab = "income(xi)", ylab = "WLS_Res",
     pch=19)
abline(h=0, lwd=2, col="red")
```

Another method of correcting for heteroscedasticity is based on the assumption that the variance of the disturbances is positively associated with level of income X , that is,

$$\sigma_i^2 = \sigma^2 X_i^2.$$

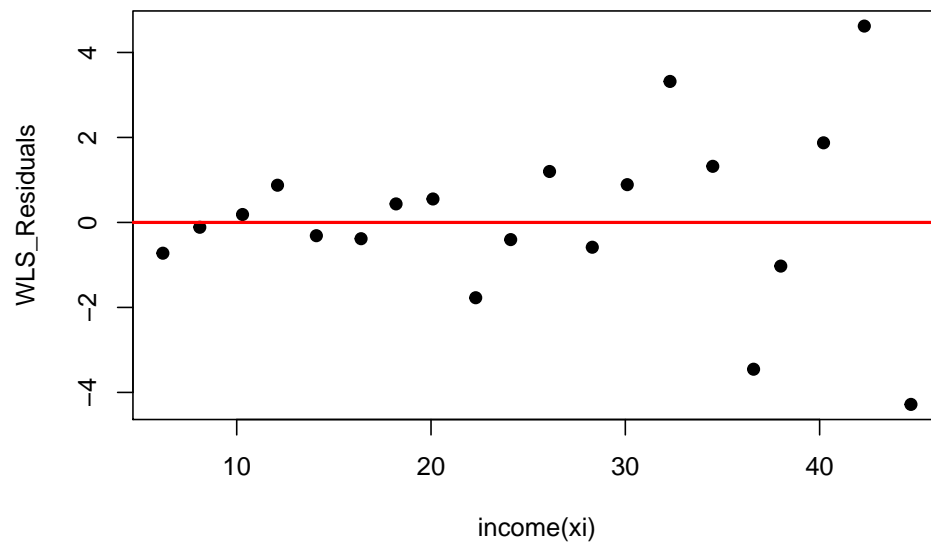


Figure 2.10: Residual Plot of WLS

The model we are going to estimate is then:

$$\frac{y_i}{x_i} = \alpha \left(\frac{1}{x_i} \right) + \beta \left(\frac{X_i}{x_i} \right) + \frac{\mathcal{E}_i}{x_i}$$

$$\Rightarrow \frac{y_i}{x_i} = \alpha \left(\frac{1}{x_i} \right) + \beta + U_i, \quad \dots \text{where } U_i = \frac{\mathcal{E}_i}{x_i}.$$

This simply means that we apply OLS by regressing $\frac{y_i}{x_i}$ on $\frac{1}{x_i}$. **(Exercise)**

2.4 Autocorrelation

objective of the lesson:

- To identify the nature of autocorrelation
- To detect the consequence of serial-autocorrelation
- To take a remedial measure for serial-autocorrelation problem
- checking spatial autocorrelation in OLS residuals
- To distinguish spatial error and lag model

Introduction

In cross-sectional studies, such as households or firms are **randomly selected** no prior reason to believe that the error term pertaining to one household or a firm is correlated with the error term of another household or firm. however, in a cross-sectional data considering a space order may exist and is called **spatial autocorrelation**. In cross-sectional analysis, the ordering of the data must have some logic, or economic interest, to make sense of any determination of whether (spatial) autocorrelation is present or not.

In time series data, where observations follow a natural ordering over time, are likely to be correlated. successive observations over short interval are also more likely to be correlated as compared to long interval difference.

2.4.1 Nature of the autocorrelation

The term autocorrelation may be defined as a correlation between order of a single variable, that is to itself. This order may across space (spatial autocorrelation) or across time(time-correlation or simply autocorrelation). One of the assumption of the CLRM is that autocorrelation does not exist in the disturbances \mathcal{E}_i .

$$E(\mathcal{E}_i \mathcal{E}_j) = 0$$

that is, the disturbance term relating to any observation is not influenced by the disturbance term relating to any other observation. For instance, in a study of the relationship between output and inputs of a firm or industry from monthly observations, non-autocorrelation of the disturbance implies that the effect of machine breakdown is strictly temporary in the sense that only the current month's output is affected. But in practice, the effect of a machine breakdown in one month may affect current month's output as well as the output of subsequent months. And, if there is such a dependence, we have autocorrelation. Mathematically,

$$E(\mathcal{E}_i \mathcal{E}_j) \neq 0 \text{ for all } i \neq j$$

In regressions involving time series data, successive observations are likely to be interdependent.

Reasons for occurring serial correlation includes:

- 1) Inertia
- 2) Specification Bias: Excluded Variables Case.
- 3) Specification Bias: Incorrect Functional Form.
- 4) Cobweb Phenomenon
- 5) Lags.
- 6) Manipulation of Data(such as smothing, interpolation or extrapolation)
- 7) Data Transformation
- 8) Nonstationarity

2.4.2 OLS estimation in the presence of serial autocorrelation

What happens to the OLS estimators and their variances in the presence of autocorrelation?

Consider two-variable regression model: $Y_t = \beta_1 + \beta_2 X_t + \mathcal{E}_t$ where $t = \text{time}$

In order to see the consequences of autocorrelation, we have to specify the nature (mathematical form) of auto-correlation or the **mechanism that generates** it, i.e., $E(\mathcal{E}_t \mathcal{E}_{t+s}) \neq 0$ (for $s \neq 0$).

As first approximation or usually we assume that the errors (disturbances) follow the first-order autoregressive scheme (abbreviated as AR(1)) or generated by the mechanism:

The error process:

$$\mathcal{E}_t = \rho \mathcal{E}_{t-1} + u_t \dots \dots \dots (2.4.1) \quad (2.40)$$

where $\rho(\text{rho})$ is known as the **coefficient of autocovariance** and where u_t is the stochastic disturbance term (also called a white noise error term) such that it satisfied the standard OLS assumptions, namely

$$E(u_t) = 0, \quad \text{var}(u_t) = \sigma_u^2, \quad \text{cov}(u_t, u_{t+s}) = 0 \text{ for } s \neq 0$$

- Equation (2.40) says that the value of the disturbance term in period t is equal to ρ times its value in the previous period plus a purely random error term.
- Equation (2.40) is also called an autoregressive process of order one (or **AR(1)**).
- The name autoregressive is appropriate because equation (2.40) can be interpreted as the regression of \mathcal{E}_t on itself lagged one period.
- It is first order because \mathcal{E}_t and its immediate past value are involved; that is, the maximum lag is 1. If the model were $\mathcal{E}_t = \rho_1 \mathcal{E}_{t-1} + \rho_2 \mathcal{E}_{t-2} + u_t$, it would be an AR(2), or second-order, autoregressive scheme, and so on.

From equation (2.40) we have (beginning from \mathcal{E}_0):

$$\begin{aligned}
\mathcal{E}_1 &= \rho \mathcal{E}_0 + u_1 \\
\mathcal{E}_2 &= \rho \mathcal{E}_1 + u_2 \\
\mathcal{E}_2 &= \rho(\rho \mathcal{E}_0 + u_1) + u_2 \\
\mathcal{E}_2 &= \rho^2 \mathcal{E}_0 + \rho u_1 + u_2 \\
\mathcal{E}_3 &= \rho \mathcal{E}_2 + u_3 \\
\mathcal{E}_3 &= \rho(\rho^2 \mathcal{E}_0 + \rho u_1 + u_2) + u_3 \\
\mathcal{E}_3 &= \rho^3 \mathcal{E}_0 + \rho^2 u_1 + \rho u_2 + u_3 \\
\mathcal{E}_t &= \rho^t \mathcal{E}_0 + \rho^{t-1} u_1 + \rho^{t-2} u_2 + \dots + \rho u_{t-1} + u_t \\
\mathcal{E}_t &= \rho^t \mathcal{E}_0 + \sum_{j=0}^{t-1} \rho^j u_{t-j}
\end{aligned}$$

For this process to be stable, ρ must be less than one in absolute value, that is, $|\rho| < 1$. In such cases, as $t \rightarrow \infty$, we have $\rho^t \mathcal{E}_0 \rightarrow 0$. Thus, \mathcal{E}_t can be expressed as:

$$\mathcal{E}_t = \sum_{j=0}^{t-1} \rho^j u_{t-j} = u_t + \rho^1 u_{t-1} + \rho^2 u_{t-2} + \rho^3 u_{t-3} + \dots$$

Since $|\rho| < 1$, ρ^j keeps on decreasing as j keeps on increasing. This means that the effect of the recent past is significant and the effect keeps on diminishing the further back we go (called **fading memory**).

The mean of \mathcal{E} is:

$$E(\mathcal{E}_t) = \sum_{j=0}^{\infty} \rho^j E(u_{t-j}) = 0$$

The variance of \mathcal{E}_t is:

$$\begin{aligned} Var(\mathcal{E}_t) &= E(\mathcal{E}_t - E(\mathcal{E}_t))^2 \\ &= E(\mathcal{E}_t^2) \\ &= E\left(\sum_{j=0}^{\infty} \rho^j u_{t-j}\right)^2 \\ &= E\left(\sum_{j=0}^{\infty} \rho^2 j u_{t-j}^2 + \sum_{j=0}^{\infty} \sum_{i \neq 0}^{\infty} \rho^j \rho^i u_{t-j} u_{t-i}\right) \\ &= \left(\sum_{j=0}^{\infty} \rho^2 j E(u_{t-j}^2) + \sum_{j=0}^{\infty} \sum_{i \neq 0}^{\infty} \rho^j \rho^i E(u_{t-j} u_{t-i})\right) \end{aligned}$$

where $E(u_{t-j}^2) = \sigma_u^2$ and $E(u_{t-j} u_{t-i}) = 0$ for $i \neq j$

$$\begin{aligned} &= \sum_{j=0}^{\infty} \rho^2 j \sigma_u^2 = \sigma_u^2 \sum_{j=0}^{\infty} \rho^2 j = \frac{\sigma_u^2}{1 - \rho^2} \\ \implies Var(\mathcal{E}_t) &= \frac{\sigma_u^2}{1 - \rho^2} = \sigma^2 \end{aligned}$$

(Assumption of constant error variance)

The covariance between \mathcal{E}_t and \mathcal{E}_{t-1} is derived as:

$$\begin{aligned}
cov(\mathcal{E}_t, \mathcal{E}_{t-1}) &= E(\mathcal{E}_t \mathcal{E}_{t-1}) \\
&= E[(u_t + \rho^1 + \rho^2 u_{t-2} + \dots)(u_{t-1} + \rho^1 u_{t-2} + \rho^2 u_{t-3} + \dots)] \\
&= E(u_t u_{t-1}) + \rho E(u_t u_{t-2}) + \rho^2 E(u_t u_{t-3}) + \dots \\
&\quad + \rho E(u_{t-1} u_{t-1})^* + \rho^2 E(u_{t-1} u_{t-2}) + \rho^3 E(u_{t-1} u_{t-3}) + \dots \\
&\quad + \rho^2 E(u_{t-2} u_{t-1}) + \rho^3 E(u_{t-2} u_{t-2})^* + \rho^4 E(u_{t-2} u_{t-3}) + \dots \\
&\quad + \rho^3 E(u_{t-3} u_{t-1}) + \rho^4 E(u_{t-3} u_{t-2}) + \rho^5 E(u_{t-3} u_{t-3})^* + \dots \\
&= \rho \sigma_u^2 + \rho^3 \sigma_u^2 + \rho^5 \sigma_u^2 + \dots \\
&= \rho \sigma_u^2 (1 + \rho^2 + \rho^4 + \dots) \\
cov(\mathcal{E}_t, \mathcal{E}_{t-1}) &= \rho \sigma_u^2 \sum_{j=0}^{\infty} (\rho^2)^j = \frac{\rho \sigma_u^2}{1 - \rho^2} = \rho \sigma^2
\end{aligned} \tag{2.41}$$

Similarly, it can be shown that:

$$\begin{aligned}
cov(\mathcal{E}_t, \mathcal{E}_{t-2}) &= E(\mathcal{E}_t, \mathcal{E}_{t-2}) = \rho^2 \sigma^2. \\
cov(\mathcal{E}_t, \mathcal{E}_{t-3}) &= E(\mathcal{E}_t, \mathcal{E}_{t-3}) = \rho^3 \sigma^2.
\end{aligned}$$

Or in general:

$$cov(\mathcal{E}_t, \mathcal{E}_{t-s}) = E(\mathcal{E}_t, \mathcal{E}_{t-s}) = \rho^s \sigma^2$$

Equivalently, we have:

$$cov(\mathcal{E}_t, \mathcal{E}_s) = E(\mathcal{E}_t, \mathcal{E}_s) = \rho^{|s-t|} \sigma^2$$

Thus, the relationship between the disturbances depends on the value of the parameter ρ . From

equation (2.41) we have:

$$\begin{aligned} cov(\mathcal{E}_t, \mathcal{E}_{t-1}) &= E(\mathcal{E}_t, \mathcal{E}_{t-1}) = \rho\sigma^2 \\ \Rightarrow \rho &= \frac{cov(\mathcal{E}_t, \mathcal{E}_{t-1})}{\sigma^2} \\ \Rightarrow \rho &= \frac{cov(\mathcal{E}_t, \mathcal{E}_{t-1})}{\sqrt{var(\mathcal{E}_t)}\sqrt{var(\mathcal{E}_{t-1})}} \end{aligned}$$

- $\Rightarrow \sigma$ is nothing but the **coefficient of correlation** between \mathcal{E}_t and \mathcal{E}_{t-1} .
- ρ is estimated from the residuals of the **OLS** regression by:

$$\hat{\rho} = \frac{\sum_{t=2}^T \hat{\mathcal{E}}_t \hat{\mathcal{E}}_{t-1}}{\sqrt{\sum_{t=1}^T \hat{\mathcal{E}}_t^2} \sqrt{\sum_{t=2}^T \hat{\mathcal{E}}_{t-1}^2}}$$

2.4.2.1 Properties of OLS estimators under serial autocorrelation The model (in deviations form) is:

$$\begin{aligned} y_t &= \beta X_t + \mathcal{E}_t \\ \mathcal{E}_t &= \rho \mathcal{E}_{t-1} + u_t, |\rho| < 1 \end{aligned}$$

where $U_t = \mathcal{E}_t - \rho \mathcal{E}_{t-1}$ satisfies all assumptions of the **CLRM** ($E(u_t) = 0, Var(u_t) = E(u_t^2) = \sigma_u^2$ and $E(u_t u_s) = 0$ for $t \neq s$).

The **OLS** estimator of β is:

$$\hat{\beta} = \frac{\sum x_t y_t}{\sum x_t^2}$$

Then we have:

$$\begin{aligned}\hat{\beta} &= \frac{\sum x_t \beta X_t + \mathcal{E}_t}{\sum x_t^2} = \beta + \frac{\sum x_t \mathcal{E}_t}{\sum x_t^2} \\ \Rightarrow E(\hat{\beta}) &= \beta + \frac{\sum x_t E(\mathcal{E}_t)}{\sum x_t^2} = \beta, \text{ where } E(\mathcal{E}_t) = 0.\end{aligned}$$

Thus, $\hat{\beta}$ is still an unbiased estimator of β .

Recall that the variance of the **OLS** estimator $\hat{\beta}$ when there is no autocorrelation is given by:

$$Var(\hat{\beta}_{OLS}) = \frac{\sigma^2}{\sum x_t^2} \quad (2.42)$$

Now under the **AR(1)** scheme, it can be shown that the variance of this estimator is:

$$\begin{aligned}Var(\hat{\beta})_{AR(1)} &= E(\hat{\beta} - \beta)^2 = E\left(\frac{\sum x_t \mathcal{E}_t}{\sum x_t^2}\right)^2 \\ &= \frac{1}{(\sum x_t^2)^2} E\left(\sum x_t^2 \mathcal{E}_t^2 + 2 \sum_{s>t} x_t \mathcal{E}_t x_s \mathcal{E}_s\right) \\ &= \frac{1}{(\sum x_t^2)^2} \sum x_t^2 E(\mathcal{E}_t^2) + \frac{2}{(\sum x_t^2)^2} \sum_{s>t} x_t x_s E(\mathcal{E}_t \mathcal{E}_s) \dots (but, E(\mathcal{E}_t \mathcal{E}_s) = \rho^{|s-t|} \sigma^2) \\ &= \frac{\sigma^2}{\sum x_t^2} + \frac{2\sigma^2}{\sum x_t^2} \left[\rho \frac{\sum x_t x_{t-1}}{\sum x_t^2} + \rho^2 \frac{\sum x_t x_{t-2}}{\sum x_t^2} + \rho^3 \frac{\sum x_t x_{t-3}}{\sum x_t^2} + \dots \right] \\ &= Var(\hat{\beta})_{OLS} + \frac{2\sigma^2}{\sum x_t^2} \left[\rho \frac{\sum x_t x_{t+1}}{\sum x_t^2} + \rho^2 \frac{\sum x_t x_{t+2}}{\sum x_t^2} + \rho^3 \frac{\sum x_t x_{t+3}}{\sum x_t^2} + \dots \right] \quad (2.43)\end{aligned}$$

Therefore, when $\rho > 0$ and x_t is positively correlated with x_{t+1}, x_{t+2}, \dots , the second term on the right hand side will usually be positive, and we have:

$$Var(\hat{\beta})_{AR(1)} = Var(\hat{\beta})_{OLS}$$

Note that the stated conditions concerning ρ and the X_t 's are fairly common in economic time series, that is, usually the terms such as $\sum x_t x_{t+1}$ or $\sum x_t x_{t+2}$ are expected to be positive since it is not unreasonable to assume that successive terms in economic time series data are positively correlated (consumption high in one period usually means it will be high in the next period as well).

Thus, if the errors are autocorrelated, and yet we persist in using **OLS**, then the variances of regression coefficients will be **under-estimated** leading to narrower confidence intervals, high values of R^2 and inflated t-ratios.

2.4.2.2 The BLUE estimator in the presence of serial-autocorrelation In the case of autocorrelation can we find an estimator that is BLUE? Yes.

Continuing with the two-variable model and assuming the AR(1) process, we can show that the BLUE estimator of β_2 is given by the following expression.

$$\hat{\beta}_2^{GLS} = \frac{\sum_{t=2}^T (x_t - \rho x_{t-1})(y_t - \rho y_{t-1})}{\sum_{t=2}^T (x_t - \rho x_{t-1})^2} + C \quad (2.44)$$

where C is a correction factor that may be disregarded in practice. And its variance is given by

$$var(\hat{\beta}_2^{GLS}) = \frac{\sigma^2}{\sum_{t=2}^T (x_t - \rho x_{t-1})^2} + C \quad (2.45)$$

where D too is a correction factor that may also be disregarded in practice.

The estimator β_2^{GLS} , as the superscript suggests, is obtained by the **method of GLS**. in GLS we incorporate any additional information we have (e.g., the nature of the heteroscedasticity or of the autocorrelation) directly into the estimating procedure by transforming the variables, whereas in OLS such side information **is not directly** taken into consideration.

Hence, under autocorrelation, it is the GLS estimator given in (2.44) that is BLUE, and the minimum variance is now given by (2.45) and not by (2.43) and obviously not by (2.42) .

2.4.2.3 Consequences of Using OLS in The Presence of Serial Autocorrelation

- 1) The residual variance $\hat{\sigma}^2 = \frac{\hat{\mathcal{E}}_i}{n-2}$ is likely to underestimate the true σ^2 .
- 2) As a result, we are likely to overestimate R^2 .
- 3) **OLS** estimators are still unbiased.
- 4) **OLS** estimators are consistent, i.e., their variances approach to zero, as the sample size gets

larger and larger.

- 5) **OLS** estimators are no longer efficient.
- 6) The estimated variances of the **OLS** estimators are biased, and as a consequence, the conventional confidence intervals and tests of significance are not valid.

2.4.2.4 Detecting Serial-Autocorrelation 1. Graphical method

Plot the estimated residuals $\hat{\mathcal{E}}_t = y_t - \hat{y}_t$ against time. If we see a clustering of neighbouring residuals on one or the other side of the line $\mathcal{E} = 0$ then such clustering is a sign that the errors are autocorrelated.

Example

Using `auto.csv` data where the variables are about investment and value of outstanding shares for the years 1935-1953. Using graphical method check that for presence of autocorrelation?

```
rm(list=ls())

auto=read.csv("auto.csv",header=T,sep=",")

fit<-lm(invest~vsh,data=auto)
plot(auto$year,fit$residuals,xlab = "Year",ylab = "Unstandrdized Residuals",
abline(h=0,lwd=2,col="red")
s<-seq(length(auto$year)-1)
arrows(auto$year[s], fit$residuals[s], auto$year[s+1], fit$residuals[s+1],
col = 1:3,lwd = 1)
```

We can see a clustering of neighbouring residuals on one or the other side of the line $\hat{\mathcal{E}}_i = 0$ (see Figure 2.11). This might be a sign that the errors are autocorrelated. However, we do not make a final judgment until we apply formal tests of autocorrelation

2. Durbin-Watson (DW) test

The **DW** test statistic is computed as:

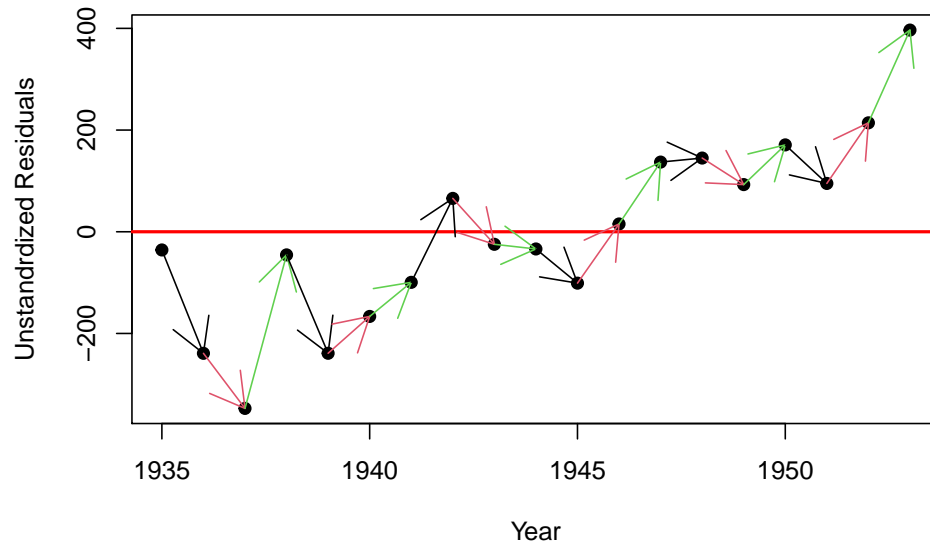


Figure 2.11: Plot of serial correlation

$$d = \frac{\sum_{t=2}^T (\hat{\mathcal{E}}_t - \hat{\mathcal{E}}_{t-1})^2}{\sum_{t=2}^T \hat{\mathcal{E}}_t^2}$$

To test of $H_0 : \rho = 0$ versus $H_A : \rho > 0$, we can use the Durbin-Watson lower (d_L) and upper (d_U) bounds (critical values).

Decision rule:

Reject H_0 if $d < d_L$

Do not reject H_0 if $d > d_U$

The test is inconclusive if $d_L < d < d_U$

Limitations of the DW test:

- a) There are certain regions where the test is inconclusive.
- b) The test is valid only when there is an intercept term in the model.
- c) The test is invalid when lagged values of the dependent variable appear as regressors.
- d) The test is valid for the **AR(1)** error scheme only.

Example

In the data `auto.csv`, check for presence of autocorrelation using DW test?

```
rm(list=ls())
auto=read.csv("auto.csv",header=T,sep=",")

dw_test<-function(x=data.frame(),dw_lower=1.180,dw_upper= 1.401){
  fit<-lm(invest~vsh,data=x)
  rsdl<-resid(fit)
  rsdl_t=rsdl[2:19]; rsdl_t_1=rsdl[1:18] # rsdl_t=rsdl[-1];rsdl_t_1=rsdl[-19]
  dw_cal<-sum((rsdl[-1]-rsdl[-19])^2)/sum(rsdl[-1]^2)

  # I recommend to use the 5 commented lines instead of the next 4
  # if(dw_cal< dw_lower){
  #   cat("Reject Ho: rho=0 ", " ,since the calculated dw-test= ",dw_cal,"<",
  #       dw_lower,"=dw_lower and concludes \n",
  #       "that a significant presence of serial autocorrelation.")
  # }

  if(dw_cal< dw_lower){
    cat("Reject Ho: rho=0 ", " ,since dw_cal= ",dw_cal,"<",
        dw_lower,"=dw_lower")
  }

  if(dw_cal>dw_lower & dw_cal<dw_upper){
    cat("The test is inconclusive")
  }

  if(dw_cal> dw_upper){
    cat("Do not reject Ho: rho=0")
  }
}
```



```
}
dw_test(x=auto)

## Reject Ho: rho=0 , since dw_cal= 0.554 < 1.18 =dw_lower
```

3. Breusch-Godfrey (BG) Test

Assume that the error term follows the **autoregressive scheme of order p (AR(p))** given by:

$$\mathcal{E}_t = \rho_1 \mathcal{E}_{t-1} + \rho_2 \mathcal{E}_{t-2} + \dots + \rho_p \mathcal{E}_{t-p} + u_t$$

where u_t fulfills all assumption of the **CLRM**. The null hypothesis to be tested is:

$$H_0 : \rho_1 = \rho_2 = \dots + \rho_p = 0$$

Steps :

1. Estimate the model:

$$Y_t = \alpha + \beta X_t + \mathcal{E}_t, \dots (t = 1, 2, \dots, T)$$

using **OLS** and obtain the residuals $\hat{\mathcal{E}}_t$.

2. Regress $\hat{\mathcal{E}}_t$ on X_t and $\hat{\mathcal{E}}_{t-1}, \hat{\mathcal{E}}_{t-2}, \dots, \hat{\mathcal{E}}_{t-p}$ that is, run the following auxiliary regression:

$$Y_t = \alpha + \beta X_t + \rho_1 \hat{\mathcal{E}}_{t-1} + \rho_2 \hat{\mathcal{E}}_{t-2} + \dots + \rho_p \hat{\mathcal{E}}_{t-p} + \xi_t$$

3. Obtain the coefficient of determination R^2 from the auxiliary regression.

4. If the sample size T is large, Breusch and Godfrey have shown that $(T - p)R^2$ follows the Chi-square (χ^2) distribution with p degrees of freedom.

Decision rule:

Reject the null hypothesis of no **AC** if $(T - p)R^2$ exceeds the tabulated value from the χ^2

distribution with p degrees of freedom for a given level of significance α .

Advantages of the BG test

- a) The test is always conclusive.
- b) The test is valid when lagged values of the dependent variable appear as regressors.
- c) The test is valid for higher order **AR** schemes (not just for **AR(1)** error scheme only).

Example

in the data `auto.csv`, check AC using BG test?

```
rm(list=ls())
auto=read.csv("auto.csv",header=T,sep=",")
T=nrow(auto)
# First apply OLS and obtain residuals
fit<-lm(invest~vsh,data=auto)
rsdl<-resid(fit)
rsdl_t=rsdl[-1];rsdl_t_1=rsdl[-19]
# to determine R2 from auxiliary regression
fit1<-lm(rsdl_t~auto$vsh[-1]+rsdl_t_1)
# (T-p)R2
(bp_cal<-(T-1)*summary(fit1)$r.squared)

## [1] 9.149

# p-value
(p_value<-pchisq(q=bp_cal, df=1,lower.tail=F))

## [1] 0.002489

#x_tab
(x_tab<-qchisq(0.05, df=1,lower.tail=F))

## [1] 3.841
```

Since the calculated test statistic(9.1485) exceeds the critical value(3.8415), we reject the null hypothesis $H_0 : \rho = 0$ and conclude that there is error AC(p-value=0.0025<0.05).

4. Test based on the partial autocorrelation function (PACF) of OLS residuals

Plot the **PACF** of **OLS** residuals. If the function at lag one is outside the 95% upper or lower confidence limits, then this is an indication that the errors follow the **AR(1)** process. Higher order error processes can be detected similarly.

Example

check using pacf from **auto.csv**

```
rm(list=ls())

rm(list=ls())
auto=read.csv("auto.csv",header=T,sep=",")
T=nrow(auto)
# First apply OLS and obtain residuals
fit<-lm(invest~vsh,data=auto)
summary(fit)

##
## Call:
## lm(formula = invest ~ vsh, data = auto)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -347.4  -100.3   -24.7   116.0   396.6
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -186.1598    216.2925  -0.86   0.4014
## vsh           0.1753     0.0497    3.53   0.0026 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Residual standard error: 185 on 17 degrees of freedom
## Multiple R-squared:  0.422, Adjusted R-squared:  0.389
## F-statistic: 12.4 on 1 and 17 DF, p-value: 0.00259
```

```
Residuals<-resid(fit)
pacf(Residuals)
```

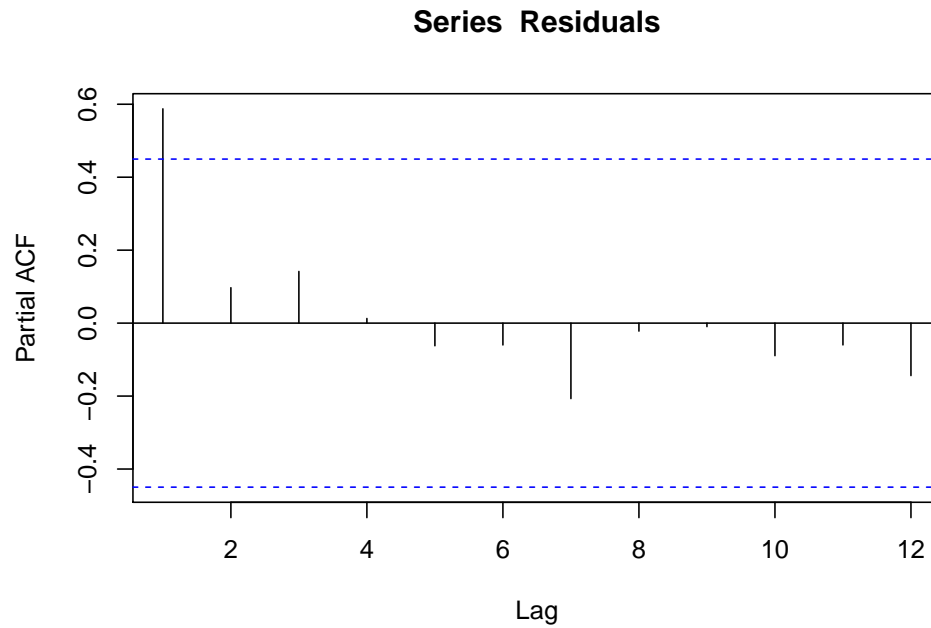


Figure 2.12: Partial autocorrelation function

The PACF at lag one is outside the 95% confidence limits and indicates that the errors follow the AR(1) process (see 2.12).

2.4.2.5 Correcting for error Autocorrelation Consider the model:

$$Y_t = \alpha + \beta X_t + \mathcal{E}_t, \dots (t = 1, 2, \dots, T) \quad (2.46)$$

where the errors are generated according to the **AR**(1) scheme:

$$\mathcal{E}_t = \rho \mathcal{E}_{t-1} + u_t, |\rho| < 1$$

Here, $U_t = \mathcal{E}_t - \rho\mathcal{E}_{t-1}$ satisfies all assumptions of the **CLRM** (that is, $E(u_t) = 0$, $var(u_t) = E(u_t^2) = \sigma_u^2$ and $E(u_t u_s) = 0$ for $t \neq s$).

Suppose by applying any one of the above tests you come to the conclusion that the errors are autocorrelated. What to do next?

Lagging equation (2.46) by one period and multiplying throughout by ρ , we get:

$$\rho Y_{t-1} = \rho\alpha + \rho\beta X_{t-1} + \rho\mathcal{E}_{t-1} \quad (2.47)$$

Subtracting equation (2.46) from equation (2.47), we get:

$$\begin{aligned} Y_t - \rho Y_{t-1} &= \alpha(1 - \rho) + \beta(X_t - \rho X_{t-1}) + (\mathcal{E}_t - \rho\mathcal{E}_{t-1}) \\ \text{let, } Y_t - \rho Y_{t-1} &= y_t^*, \alpha(1 - \rho) = \alpha^*, \beta(X_t - \rho X_{t-1}) = x_t^* \text{ and } (\mathcal{E}_t - \rho\mathcal{E}_{t-1}) = U_t \\ &\Rightarrow Y_t^* = \alpha^* + x_t^* + U_t \end{aligned} \quad (2.48)$$

The above transformation is known as the **Cochrane-Orcutt transformation**. Since

$U_t = \mathcal{E}_t - \rho\mathcal{E}_{t-1}$ fulfils all assumption of the **CLRM**, we can apply **OLS** to (2.48) to get estimators which are **BLUE**.

Problem : The above transformation requires a knowledge of the value of ρ . Thus we need to estimate it.

Methods of estimation of ρ

a) Using the Durbin-Watson statistic.

It can be shown that as T (the sample size) gets larger, the DW statistic d approaches to $2(1 - \rho)$, i.e., $d \rightarrow 2(1 - \rho)$ as $T \rightarrow \infty$. Thus, we can use this fact to construct an estimator of ρ as:

$$\hat{\rho} = 1 - \frac{d}{2}$$

Note : This estimator is **highly inaccurate** if the sample size is **small**.

```
# from previous dw=
d=0.5539826
(rho=1-d/2)
```

```
## [1] 0.723
```

and hence the estimated $\hat{\rho} = 0.723$.

b) From **OLS** residuals

Regress OLS residuals $\hat{\mathcal{E}}_t$ on $\hat{\mathcal{E}}_{t-1}$ **with out a constant term**:

$$\hat{\mathcal{E}}_t = \rho \hat{\mathcal{E}}_{t-1} + u_t$$

An estimate of ρ is the estimated coefficient of $\hat{\mathcal{E}}_{t-1}$.

```
rm(list=ls())
auto=read.csv("auto.csv",header=T,sep=",")
fit<-lm(invest~vsh,data=auto)
rsdl<-resid(fit)
fit<-lm(rsdl[-1]~rsdl[-19]-1)
summary(fit)$coeff
```

```
##              Estimate Std. Error t value Pr(>|t|)
## rsdl[-19]    0.8049      0.2057    3.913  0.00112
```

and hence the estimated $\hat{\rho} = 0.8049$.

c) Durbin's method

Run the regression of Y_t on Y_{t-1} , X_t and X_{t-1} :

$$y_t = \delta + \rho Y_{t-1} + \beta x_t + \beta \rho x_{t-1} + u_t$$

An estimator of ρ is the estimated coefficient of Y_{t-1} .

```
rm(list=ls())
auto=read.csv("auto.csv",header=T,sep=",")
```

```

y=auto$invest
x=auto$vsh
y_t1=y[-1]; y_t=y[-19]
x_t1=x[-1]; x_t=x[-19]
fit<-lm(y_t~y_t1+x_t+x_t1)
summary(fit)$coeff

```

```

##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  264.28692    87.60698   3.017 9.240e-03
## y_t1         0.72866     0.06566  11.098 2.534e-08
## x_t          0.04010     0.01615   2.483 2.631e-02
## x_t1        -0.07608     0.01761  -4.320 7.063e-04

```

and hence the estimated $\hat{\rho} = 0.7287$.

All tests indicated that there is serial autocorrelation. Thus, we need to apply the Cochrane-Orcutt transformation. To obtain an estimate $\hat{\rho}$ of ρ , let's use the result of regressing the OLS residuals $\hat{\mathcal{E}}_t$ on \mathcal{E}_{t-1} with out a constant term (as in b above) and is $\hat{\rho} = 0.8049$.

Then apply the following (Cochrane-Orcutt) transformation:

$$\begin{aligned}
 Y_t - \rho Y_{t-1} &= \alpha(1 - \rho) + \beta(X_t - \rho X_{t-1}) + u_t \\
 \Rightarrow Y_t^* &= \alpha^* + X_t^* + u_t
 \end{aligned}
 \tag{2.49}$$

Note that equation (2.49) fulfills all basic assumptions and, thus, we can estimate the parameters in this equation by an OLS procedure. Using $\hat{\rho} = 0.8049$, we obtain Y_t^* (invst_trnsf) and X_t^* (vsh_trnsf) and estimate the regression of Y_t^* on X_t^* . The results are:

```

rm(list=ls())
auto=read.csv("auto.csv", header=T, sep=",")
# lets use estiamte of rho from OLS in (b)
rho=0.8049

```

```

y=auto$invest
x=auto$vsh
y_t1=y[-1];y_t=y[-19]
invst_trnsf=y_t-rho*y_t1
x_t1=x[-1]; x_t=x[-19]
vsh_trnsf=x_t-rho*x_t1

fit_trnsf<-lm(invst_trnsf~vsh_trnsf)
summary(fit_trnsf)

##
## Call:
## lm(formula = invst_trnsf ~ vsh_trnsf)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -163.4   -28.2    18.9    36.7    77.2
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   11.4528    18.3202   0.63  0.54069
## vsh_trnsf      0.0686     0.0160   4.29  0.00056 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 63.1 on 16 degrees of freedom
## Multiple R-squared:  0.535, Adjusted R-squared:  0.506
## F-statistic: 18.4 on 1 and 16 DF, p-value: 0.000557

```

The partial autocorrelation function of the residuals in the transformed model is shown below. It can be seen that the function lies within the upper and lower confidence limits, indicating that the

autocorrelation structure has been properly dealt with (see Figure 2.13).

```
Residuals<-resid(fit_trnsf)
pacf(Residuals)
```

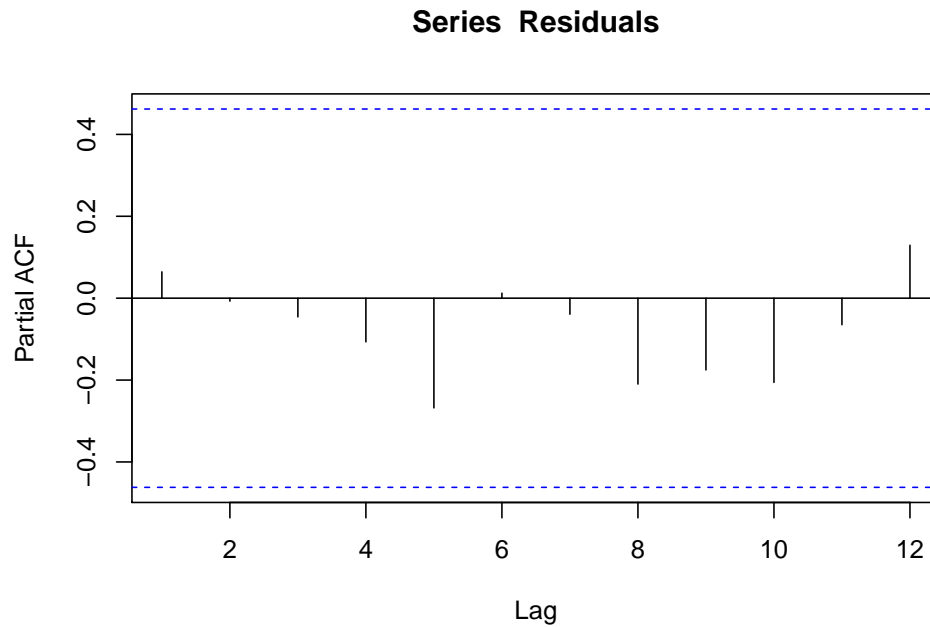


Figure 2.13: Partial autocorrelation function for Cochrane-Orcutt transformation

Generalized Least Squares (GLS)

Introduction

Consider the model:

$$y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k x_{ki} + \mathcal{E}_i, \dots, i = 1, 2, \dots, n$$

In matrix form, the above model is:

$$Y = X\beta + \varepsilon$$

where:

$$Y_{(nx1)} = \begin{bmatrix} Y_1 \\ Y_2 \\ \cdot \\ \cdot \\ \cdot \\ Y_n \end{bmatrix} \quad x_{(nxk)} = \begin{bmatrix} 1 & X_{21} & \cdot & \cdot & \cdot & X_{k1} \\ 1 & X_{22} & \cdot & \cdot & \cdot & X_{k2} \\ \cdot & \cdot & & & & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & & & & \cdot \\ 1 & X_{2n} & \cdot & \cdot & \cdot & X_{kn} \end{bmatrix} \quad \beta_{(kx1)} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \cdot \\ \cdot \\ \cdot \\ \beta_k \end{bmatrix} \quad \text{and} \quad \mathcal{E}_{(nx1)} = \begin{bmatrix} \mathcal{E}_1 \\ \mathcal{E}_2 \\ \cdot \\ \cdot \\ \cdot \\ \mathcal{E}_n \end{bmatrix}$$

Two of the assumptions of the CLRM are:

- a) $E(\mathcal{E}_i \mathcal{E}_j) = 0$ for $i \neq j$ (no autocorrelation)
- b) $E(\mathcal{E}_i^2) = \sigma^2$ for all i (no heteroscedasticity).

When these assumptions are fulfilled, the variance-covariance matrix of the error vector \mathcal{E} is given by :

$$\sum_{\mathcal{E}} = E \{ \mathcal{E}(\mathcal{E})' \} = \sigma^2 I_n$$

where I_n is an identity matrix of size n by n .

However, these two assumptions may not hold in real life and actual data we take for analysis may not fulfill them. When these assumptions do not hold, OLS is no more efficient. Instead we use another method of estimation, called GLS, which can solve these problems.

Note that these assumptions are violated if either:

- a) the errors are heteroscedastic: $Var(\mathcal{E}_i) = E(\mathcal{E}_i^2) = \sigma_i^2$, or,
- b) the errors are autocorrelated: $Cov(\mathcal{E}_i, \mathcal{E}_j) = E(\mathcal{E}_i \mathcal{E}_j) = \sigma_{ij} \neq 0$ for $i \neq j$

In such cases, the variance-covariance matrix of the error vector \mathcal{E} is given by:

$$\sum_{\mathcal{E}} = E\{\mathcal{E}\mathcal{E}'\} = \begin{bmatrix} E(\mathcal{E}_1^2) & E(\mathcal{E}_1\mathcal{E}_2) & \dots & E(\mathcal{E}_1\mathcal{E}_n) \\ E(\mathcal{E}_2\mathcal{E}_1) & E(\mathcal{E}_2^2) & \dots & E(\mathcal{E}_2\mathcal{E}_n) \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & & \cdot \\ E(\mathcal{E}_n\mathcal{E}_1) & E(\mathcal{E}_n\mathcal{E}_2) & \dots & E(\mathcal{E}_n^2) \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2n} \\ \cdot & \cdot & & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & & \cdot \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{bmatrix} = \sigma^2\Omega \neq \sigma^2I_n$$

The GLS estimation of β

Generally, given the model: $Y = X\beta + \mathcal{E}$, we use GLS to estimate the coefficient vector β when the assumption that:

$$Cov(\mathcal{E}) = E(\mathcal{E}\mathcal{E}') = \sigma^2I_n$$

is not fulfilled, and instead we have:

$$Cov(\mathcal{E}) = E(\mathcal{E}\mathcal{E}') = \sigma^2\Omega \neq \sigma^2I_n$$

where Ω is symmetric ($\Omega = \Omega'$) and positive definite ($\Omega > 0$).

The spectral decomposition of Ω is given by:

$$\Omega = B\Lambda B'$$

where Λ is a diagonal matrix whose diagonal elements are the eigenvalues of Ω (all of which are positive since Ω is positive definite) and B is a matrix whose columns are the eigenvectors corresponding to the eigenvalues of Ω .

Since all eigenvalues are positive, we can write Λ as:

$$\Lambda = \Lambda^{\frac{1}{2}}\Lambda^{\frac{1}{2}}$$

where $\Lambda^{\frac{1}{2}}$ is a diagonal matrix whose elements are the square roots of the eigenvalues of Ω .

Then we have:

$$\Omega = B\Lambda B' = B\Lambda^{\frac{1}{2}}\Lambda^{\frac{1}{2}}B' = (B\Lambda^{\frac{1}{2}})(B\Lambda^{\frac{1}{2}})' = RR' \quad (2.50)$$

where $R = B\Lambda^{\frac{1}{2}}$. Then it follows that:

$$\Omega^{-1} = (RR')^{-1} = (R^{-1})'(R^{-1}) = P'P \quad (2.51)$$

where $P = R^{-1}$. Note that:

$$(PR)(PR)' = PRP'R' = R^{-1}RR'(R^{-1})' = R^{-1}R(R^{-1}R)' = I \quad (2.52)$$

Our model is:

$$Y = X\beta + \mathcal{E} \quad (2.53)$$

with $Cov(\mathcal{E}) = E\{\mathcal{E}\mathcal{E}'\} = \sigma^2\Omega \neq \sigma^2I_n$. Pre-multiply equation (2.50) by the matrix P:

$$\begin{aligned} \underbrace{PY} &= \underbrace{PX}\beta + \underbrace{P\mathcal{E}}, \text{ where } \underbrace{PY} = Y^*, \underbrace{PX}\beta = X^* \text{ and } \underbrace{P\mathcal{E}} = \mathcal{E}^* \\ &\Rightarrow Y^* = X^*\beta + \mathcal{E}^* \end{aligned} \quad (2.54)$$

Now, the variance-covariance matrix of the errors \mathcal{E}^* in the transformed model (2.54) is:

$$\begin{aligned} Cov(\mathcal{E}^*) &= E(\mathcal{E}^*\mathcal{E}^{*'}) = E[(P\mathcal{E})(P\mathcal{E}')] \\ &= P[E(\mathcal{E}\mathcal{E}')]P' \\ &= P[\sigma^2\Omega]P' \end{aligned}$$

$$= \sigma^2 P\Omega P', \dots \text{but } \Omega = RR', \dots \text{from (2.50)}$$

$$= \sigma^2 P(RR')P'$$

$$= \sigma^2(PR)(PR)' \dots (but..(PR)(PR)' = I, \dots \text{from (2.52)})$$

$$= \sigma^2 I$$

That is, $cov(\mathcal{E}^*) = E(\mathcal{E}^* \mathcal{E}^{*'}) = \sigma^2 I$. Thus, the transformed model (2.54) fulfills all assumptions of the classical linear regression model. This means that we can apply ordinary least squares (OLS) to model (2.54) to get the BLUE $\tilde{\beta}$ of β . The BLUE β is:

$$\begin{aligned} \tilde{\beta} &= (X^{*'} X^*)^{-1} X^{*'} Y^* \\ &= [(PX)'(PX)]^{-1} [(PX)'(PY)] \\ &= [X'(P'P)X]^{-1} [X'(P'P)Y] \\ &= [X'\Omega^{-1}X]^{-1} X'\Omega^{-1}Y \end{aligned}$$

This estimator $\tilde{\beta} = [X'\Omega^{-1}X]^{-1} X'\Omega^{-1}Y$ is called the **generalized least squares (GLS)** or **Aitken** estimator of β .

Estimation of $Cov(\tilde{\beta})$

$$\begin{aligned} \tilde{\beta} &= [X'\Omega^{-1}X]^{-1} X'\Omega^{-1}Y \\ &= [X'\Omega^{-1}X]^{-1} X'\Omega^{-1}(X\beta + \mathcal{E}) \\ &= [X'\Omega^{-1}X]^{-1} (X'\Omega^{-1}X)\beta + [X'\Omega^{-1}X]^{-1} X'\Omega^{-1}\mathcal{E} \\ &= \beta + [X'\Omega^{-1}X]^{-1} X'\Omega^{-1}\mathcal{E} \end{aligned}$$

The expected value of $\tilde{\beta}$ is:

$$\begin{aligned} E(\tilde{\beta}) &= E(\beta + [X'\Omega^{-1}X]^{-1} X'\Omega^{-1}\mathcal{E}) \\ &= \beta + [X'\Omega^{-1}X]^{-1} X'\Omega^{-1} \underbrace{E(\mathcal{E})}_{=0} = \beta, \dots \text{where, } E(\mathcal{E}) = 0. \end{aligned}$$

Thus, $\tilde{\beta}$ is an unbiased estimator of β . The variance-covariance matrix of $\tilde{\beta}$ is:

$$\begin{aligned}
Cov(\tilde{\beta}) &= E[(\beta - \tilde{\beta})(\beta - \tilde{\beta})'] \\
&= E[(X'\Omega^{-1}X)^{-1}X'\Omega^{-1}\mathcal{E})([X'\Omega^{-1}X]^{-1}X'\Omega^{-1}\mathcal{E})'] \\
&= \left\{ [(X'\Omega^{-1}X)^{-1}X'\Omega^{-1}] \underbrace{E(\mathcal{E}\mathcal{E}')} \left\{ ([X'\Omega^{-1}X]^{-1}X'\Omega^{-1})' \right\} \dots \text{let, } E(\mathcal{E}\mathcal{E}') = \right. \\
&= \left\{ [(X'\Omega^{-1}X)^{-1}X'\Omega^{-1}] \sigma^2 \Omega \left\{ ([X'\Omega^{-1}X]^{-1}X'\Omega^{-1})' \right\} \right. \\
&= \sigma^2 \left\{ [(X'\Omega^{-1}X)^{-1}X'\Omega^{-1} \underbrace{\Omega\Omega^{-1}} X [(X'\Omega^{-1}X)^{-1}] \right\} \dots \text{let, } \Omega\Omega^{-1} = I \\
&= \sigma^2 \left\{ [(X'\Omega^{-1}X)^{-1}X'\Omega^{-1} \underbrace{X'\Omega^{-1}X [(X'\Omega^{-1}X)^{-1}]} \right\} \dots
\end{aligned}$$

$$\text{let, } X'\Omega^{-1}X[(X'\Omega^{-1}X)^{-1}] = I$$

$$= \sigma^2 [X'\Omega^{-1}X]^{-1}$$

$$\text{Thus, } Cov(\tilde{\beta}) = \sigma^2 [X'\Omega^{-1}X]^{-1}$$

Two stage least squares (2-SLS)

Suppose the errors are autocorrelated and follow the AR(1) process, that is,

$$\mathcal{E}_t = \rho\mathcal{E}_{t-1} + U_t$$

where $E(u_t) = 0$, $var(u_t) = E(u_t^2) = \sigma_u^2$ and $E(u_t u_s) = 0$ for $t \neq s$.

We have seen that (refer to Unit 4):

$$\sigma^2 = Var(\mathcal{E}_t) = \frac{\sigma_u^2}{1 - \rho^2}$$

$$Cov(\mathcal{E}_t, \mathcal{E}_s) = E(\mathcal{E}_t, \mathcal{E}_s) = \rho^{|s-t|} \sigma^2$$

The variance-covariance matrix of the error vector \mathcal{E} is given by:

$$\begin{aligned}
Cov(\mathcal{E}) = E(\mathcal{E}\mathcal{E}') &= \begin{bmatrix} E(\mathcal{E}_1^2) & E(\mathcal{E}_1\mathcal{E}_2) & . & . & . & E(\mathcal{E}_1\mathcal{E}_T) \\ E(\mathcal{E}_2\mathcal{E}_1) & E(\mathcal{E}_2^2) & . & . & . & E(\mathcal{E}_2\mathcal{E}_T) \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ E(\mathcal{E}_T\mathcal{E}_1) & E(\mathcal{E}_T\mathcal{E}_2) & . & . & . & E(\mathcal{E}_T^2) \end{bmatrix} \\
&= \begin{bmatrix} \sigma^2 & \rho\sigma^2 & . & . & . & \rho^{T-1}\sigma^2 \\ \rho\sigma^2 & \sigma^2 & . & . & . & \rho^{T-2}\sigma^2 \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ \rho^{T-1}\sigma^2 & \rho^{T-2}\sigma^2 & . & . & . & \sigma^2 \end{bmatrix} = \sigma^2 \underbrace{\begin{bmatrix} 1 & \rho & . & . & . & \rho^{T-1} \\ \rho & 1 & . & . & . & \rho^{T-2} \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ \rho^{T-1} & \rho^{T-2} & . & . & . & 1 \end{bmatrix}} \\
&= let, \begin{bmatrix} 1 & \rho & . & . & . & \rho^{T-1} \\ \rho & 1 & . & . & . & \rho^{T-2} \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ \rho^{T-1} & \rho^{T-2} & . & . & . & 1 \end{bmatrix} = \Omega
\end{aligned}$$

$$Cov(\mathcal{E}) = E(\mathcal{E}\mathcal{E}') = \sigma^2\Omega$$

Note that Ω depends only on ρ . The 2-SLS procedure is:

1. Estimate ρ by:

$$\hat{\rho} = \frac{\sum_{t=2}^T \hat{\mathcal{E}}_t \hat{\mathcal{E}}_{t-1}}{\sum_{t=2}^T \hat{\mathcal{E}}_{t-1}^2}$$

where $\hat{\mathcal{E}}_t$ are OLS residuals.

2. Estimate Ω by:

$$\hat{\Omega} = \begin{bmatrix} 1 & \rho & . & . & . & \rho^{T-1} \\ \rho & 1 & . & . & . & \rho^{T-2} \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ \rho^{T-1} & \rho^{T-2} & . & . & . & 1 \end{bmatrix}$$

3. The two stage least squares (2-SLS) estimator of β is then given by:

$$\tilde{\beta}_{\hat{\Omega}} = [X' \hat{\Omega} X]^{-1} X' \hat{\Omega} Y$$

2.4.3 2.4.2 Spatial autocorrelation

“All places are related but nearby places are more related than distant places”. Spatial aut-ocorrelation is the formal property that measures the degree to which near and distant things are related. Statistical test of match between locational similarity and attribute similarity positive, negative or zero relationship.

2.4.3.1 Spatial Weights Matrices What is spatial neighbor? how do we measure it?

neighbor can be created based on these

1) Contiguity Based Neighbors: Areas sharing any boundary point (**QUEEN**) are taken as neighbors.

If contiguity is defined as areas sharing more than one boundary point (**ROOK**)

2) Distance based neighbors: this can be K nearest neighbors and based neighbors specified distance

Once our list of neighbors has been created, we assign spatial weights to each relationship

2.4.3.2 Row-standardized weights matrix Row standardization is used to create **proportional weights** in cases where features have an unequal number of neighbors.

Divide each neighbor weight for a feature by the sum of all neighbor weights

- Obs i has 3 neighbors, each has a weight of 1/3
- Obs j has 2 neighbors, each has a weight of 1/2

2.4.3.3 2.4.2.1 Exploratory spatial data analysis (ESDA) Exploratory Spatial Data Analysis is a set of techniques aimed at, describing and visualizing spatial distributions, identifying atypical localization or spatial outliers, detecting patterns of spatial association, clusters or hot spots, and suggesting spatial regimes or other forms of spatial heterogeneity.

These techniques provide measures of global and local spatial autocorrelation. Spatial autocorrelation can be defined as the coincidence of value similarity with location similarity or dissimilarity. Therefore, there is positive spatial autocorrelation when high or low values of a random variable tend to cluster in space, and there is negative spatial autocorrelation when geographical areas tend to be surrounded by neighbors with very dissimilar values. Spatial dependence (spatial autocorrelation) can be measured by global and local indicators:

2.4.3.3.1 2.4.2.1.1 Measure of Global Autocorrelation The essence of spatial autocorrelation is that of spatial dependency, the situation whereby observations drawn from different locations are not independent of each other. The spatial association of data collected in space is tested using a Global Moran's I, which measures similarities and dissimilarities in observations across space. The statistic, Moran's I, also uses to check for the existence of spatial autocorrelation among particular data vector or residuals as the two dimensional analog of a test for univariate time series correlation. The Moran's I is defined as:

$$I = \frac{n}{\sum_i \sum_j W_{ij}} \frac{\sum_i \sum_j W_{ij} (x_i - \bar{x})(x_j - \bar{x})}{\sum_i (x_i - \bar{x})^2}$$

The matrix form is $I = \frac{n}{S_0} \frac{Z'WZ}{Z'Z}$, where n is number of locations, S_0 is scaling factor equal to sum of all elements of neighbor matrix W and Z is vector of n observation deviation from mean.

For row standardized W the statistic is reduced to $I = \frac{Z'WZ}{Z'Z}$.

The value of I is between -1 and 1; $I = -1$ perfect negative spatial autocorrelation, $I = 1$ is perfect positive spatial autocorrelation, and $I = 0$ signifies no spatial correlation. Inference on Moran's I take normal assumption and randomization or permutation approaches to determine the distribution of test for spatial autocorrelation under null hypothesis. For data or residuals normally distributed under null hypothesis the values of I larger than the expected value of I indicate positive spatial autocorrelation, while values of I smaller than the expected indicate negative spatial autocorrelation. Alternatively, one can apply a randomization approach in which the observed values on the random variable are randomly permuted across allocations with the assumption of each value can be equally likely observed at each location to obtain reference distribution. Inference is based on the permutation approach carried out by permuting the observed values over all locations and by re computing I for each new sample.

2.4.3.3.2 2.4.2.1.2 Measures of Local Autocorrelation These measures are used when there is no global autocorrelation, and in case where measure of global does not enable us to appreciate the regional structure of spatial autocorrelation. One can wonder which regions contribute more to the global spatial autocorrelation, whether there are local spatial clusters of high or low values, and finally to what point the global evaluation of spatial autocorrelation masks atypical localizations or pockets of local non stationary. The analysis of local spatial autocorrelation is carried out with two tools. First, the Moran scatter plot which is used to visualize local spatial instability , and second local indicators of spatial association which are used to test the hypothesis of random distribution by comparing the values of each specific localization with the values in the neighboring localizations.

2.4.3.3.3 Diagnostics tests of Spatial Dependence in OLS Regression Residuals: Standard multiple linear (OLS) regression model with spatially autocorrelated residuals may violate the independence assumption for error term, consequently regression parameter estimate are no longer BLUE, consistency, and unbiased so statistical inference is unreliable. Hence the important issue in empirical spatial analysis is how one can detect the presence of spatial effects, and moreover, how one can distinguish between spatial dependence as a nuisance and a substantive spatial process. The following tests are applying to check the presence of spatial autocorrelation in OLS regression model residuals.

a) Moran's Test for Regression Residuals

A well known test for spatial autocorrelation in the regression residuals similar to Moran's I; the statistic is defined for existence of spatial autocorrelation with in standard linear regression residuals.

The Moran's I is given as:

$$I = \frac{n}{S_0} \frac{U'WU}{U'U}$$

where, U is the vector of residuals.

Statistical inference can be based on the assumption of asymptotic normality, but an exact approach depending on the matrix X is available too, although rather cumbersome to apply. Moran's I for regression residuals is a locally best invariant test and based on moments estimation. For normally distributed errors the distribution of the standardized Moran's statistic is shown to be asymptotically normal. In order to carry out an operational test, both the expected value and the variance of I are needed.

b) Lagrange Multiplier (LM) Test

The LM test statistic for spatial error is defined as:

$$I = \frac{1}{T} \frac{U'WU}{S^2}$$

where, S^2 is the maximum likelihood estimate of variance and T is a scalar computed as the trace of a quadratic expression in the weight matrix which is given by $T = tr(W'W + W^2)$.

To model the determinant factors related with crime distribution in Addis Ababa city, Ethiopia, we have spatial data as **crime_AA.shp**. To read the shape file, we need to have install some packages such as, **rgdal**, **maptools**, **tmap**, **spdep** and **sf**.

The variables included in the data are:

✓ **UNEM_R**: Unemployment rate; is the percentage of unemployed population divide by a total of economically active people.

✓ **P_LONEH**: Percentage of single parent households; the percentage of total single woman and man households divide by total households.

Table 2.6: Crme data of Addis Ababa city, 2012

	QUEEN1	UK_NAME	UK_CODE	KK_NAME
0	1	01/03	001	Akaki Kaliti_1
1	2	02/04	002	Akaki Kaliti_2
2	3	05/06	003	Akaki Kaliti_3
3	4	07/08/09	004	Akaki Kaliti_4
4	5	10/11	005	Akaki Kaliti_5
5	6	12/13	006	Akaki Kaliti_6
6	7	Kilinito, Feche Roye	007	Akaki Kaliti_7
7	8	Golanigora	008	Akaki Kaliti_8

✓ **P_HOMEOW**: Percentage home owner households; total home owner households divide by total households.

✓ **FHEADED_P**: Percentage of female headed household; percentage of total female headed household divide by total households

✓ **YOUNGM_P**: Percentage of young male aged 15-24; percentage of total male aged 15-24 divide by total population.

✓ **POP_D**: Population density; is the total population divide by area in square Kilometer.

✓ **EDUP_1000**: Percentage of educated people; the percent of the population aged 20 or older with at least some college education.

```
# install.packages(c("rgdal", "maptools", "sf", "tmap", "spdep"))

library(maptools)

crime_A<-readShapePoly("crime_data\\Crime_AA.shp")
```

```
## Warning: readShapePoly is deprecated; use rgdal::readOGR or sf::st_read
```

```
#head(crime_A@data[,c(1,22,26,29)],20)
```

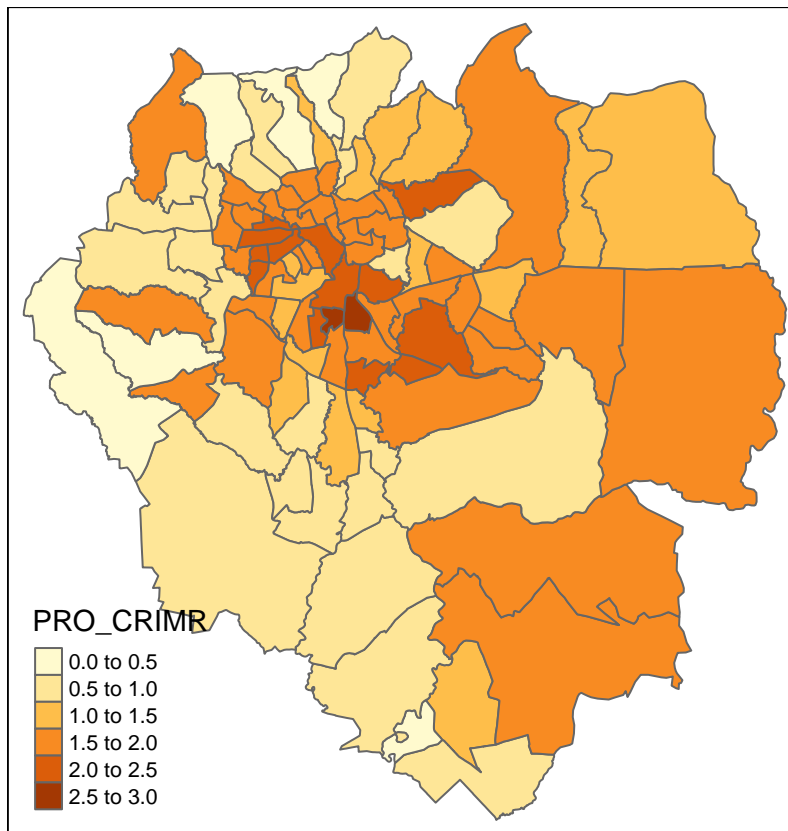
The map of crime distribution in the below shows crime are distributed highly some where and not random. However, it needs a test of spatial auto-correlation.

```
library(tmap)

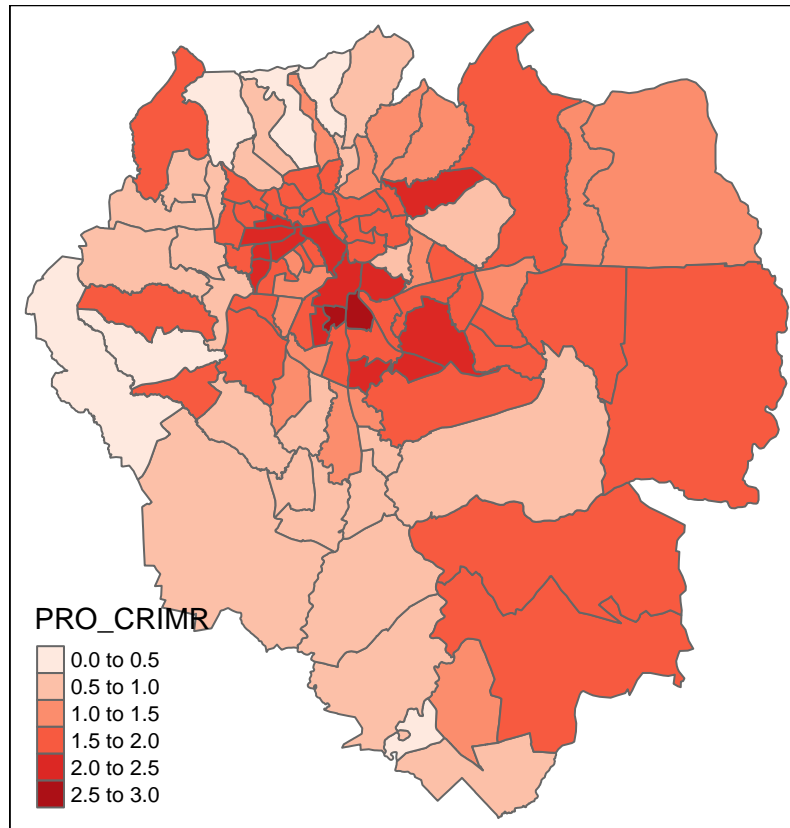
qtm(crime_A, "PRO_CRIMR") #show a variable change over space) quickly and ea
```

Table 2.7: Crme data of Addis Ababa city, 2012

	QUEEN1	UK_NAME	UK_CODE	KK_NAME
93	94	17/19/20	007	Bole_7
94	95	01	008	Bole_8
95	96	02	009	Bole_9
96	97	10	010	Bole_10
97	98	11	011	Bole_11
98	99	13/14	008	Yeka_8



```
qtm(shp=crime_A, fill="PRO_CRIMR", fill.palette="Reds") # d/t brightness of r
```



**** How to create spatial weight Matrix****

1. row standardized spatial weight m , nb2listw(neighbors info variable)
2. binary spatial weight m, nb2listw(neighbors info variable ,style="B")

row standardized spatial weight Matrix from queen

```
library(spdep)
c_nbq<-poly2nb(crime_A)

co<-coordinates(crime_A)
# To plot the neighbor coordinates
plot(crime_A)
plot(c_nbq,co,add=T)
```



```
c_nbq_w<-nb2listw(c_nbq)
# Binary spatial weight Matrix from queen
c_nbq_b<-nb2listw(c_nbq, style="B")
```

Global Spatial autocorrelation

In R, use command **moran.test()** to determine moran's index

Using `moran.test(variable,listw=)` for one sided

`moran.test(variable,listw=, alternative="two.sided")`

```
moran.test(crime_A$PRO_CRIMR, listw=c_nbq_w)
```

```
##
## Moran I test under randomisation
##
## data: crime_A$PRO_CRIMR
## weights: c_nbq_w
```

```
##
## Moran I statistic standard deviate = 8.5, p-value <2e-16
## alternative hypothesis: greater
## sample estimates:
## Moran I statistic      Expectation      Variance
##           0.505006           -0.010204           0.003635
```

the moran's I is 0.505 and indicates that there is a positive spatial autocorrelation in crime distribution($p\text{-value}=2.2e^{16} < 0.05$). This implies that OLS estimates are not BLUE and use either spatial error model or spatial lag model to determine the factors associated with crime distribution. Since, spatial statistics is delivered as one course for undergraduate statistics students, topics related to spatial error or lag model is beyond this module.

2.5 Multicollinearity

Introduction

Perfect multicollinearity means that we cannot separate the influence on Y of the independent variables that are perfectly related. assumption 6 of no perfect multicollinearity is needed to guarantee a unique solution of the OLS normal equations. Note that it applies to perfect linear relationships and does not apply to perfect non-linear relation-ships among the independent variables. In other words, one can include X_{1i} and X_{2i} like (years of experience) and (years of experience) in an equation explaining earnings of individuals. Although, there is a perfect quadratic relationship between these independent variables, this is not a perfect linear relationship and therefore, does not cause perfect multicollinearity In the construction of an econometric model, it may happen that two or more variables giving rise to the same piece of information are included, that is, we may have redundant information or unnecessarily included related variables. This is what we call a **multicollinearity(MC)** problem.

Dependent variable: Y of size $n \times 1$

Independent (explanatory) variables: X_2, X_3, \dots, X_k each of size $n \times 1$

$$\begin{aligned} \text{Model : } \mathbf{Y} &= \beta_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k + \mathcal{E} \\ &= \sum_{j=1}^k \beta_j X_j + \mathcal{E} \end{aligned}$$

where $\beta = [\beta_1, \beta_2, \beta_3, \dots, \beta_K]$, $X = [X_1, X_2, X_3, \dots, X_k]$ and X_1 is an $(n \times 1)$ column vector of ones.

One of the assumptions of the CLRM is that there is no exact linear relationship exists between any of the explanatory variables. When this assumption is violated, we speak of **perfect MC**. If all explanatory variables are uncorrelated with each other, we speak of **absence of MC**. These are two extreme cases and rarely exist in practice. Of particular interest are cases in between: **moderate to high degree of MC**.

Such kind of **MC** is so common in macroeconomic time series data (such as **GNP**, money supply, income, etc) since economic variables tend to move together over time.

2.5.1 Consequences of perfect MC

We say there is a perfect **MC** if two or more explanatory variables are perfectly correlated, that is, if the following relationship exists between the explanatory variables:

$$\alpha_2 X_2 + \alpha_3 X_3 + \dots + \alpha_k X_k = 0 \quad (2.55)$$

where $\alpha_2, \alpha_3, \dots, \alpha_k$ are constants (real numbers) such that at least two of them are not zero. One consequence of perfect **MC** is **non-identifiability** of the regression coefficient vector β . This means that one can not distinguish between two different models: $\mathbf{Y} = \mathbf{X}\beta + \mathcal{E}$ and $\mathbf{Y} = \mathbf{X}\tilde{\beta} + \mathcal{E}$. These two models are said to be **observationally equivalent**.

Illustration:

In equation (2.55), let $\alpha_2 = -1$. Then we have:

$$\begin{aligned}
X_2 &= \alpha_3 X_3 + \dots + \alpha_k X_k \\
&= \sum_{j=3}^k \alpha_j X_j
\end{aligned} \tag{2.56}$$

Then, for any scalar $\lambda \neq 0$ we have:

$$\begin{aligned}
X\beta &= \sum_{j=1}^k \beta_j X_j \\
&= \beta_1 X_1 + \beta_2 X_2 - \lambda \beta_2 X_2 + \sum_{j=3}^k \beta_j X_j + \lambda \beta_2 X_2 \\
&= \beta_1 X_1 + \beta_2 X_2 - \lambda \beta_2 X_2 + \sum_{j=3}^k \beta_j X_j + \lambda \beta_2 \sum_{j=3}^k \alpha_j X_j
\end{aligned}$$

(from equation (2.56))

$$\begin{aligned}
&= \beta_1 X_1 + (1 - \lambda) \beta_2 X_2 + \sum_{j=3}^k (\beta_j + \lambda \beta_2 \alpha_j) X_j \\
&\text{let, } \beta_1 = \tilde{\beta}_1, (1 - \lambda) \beta_2 = \tilde{\beta}_2 \text{ and } (\beta_j + \lambda \beta_2 \alpha_j) = \tilde{\beta}_j \\
&\text{then, } X\beta = \sum_{j=1}^k \tilde{\beta}_j X_j = X\tilde{\beta}
\end{aligned}$$

Thus, one can not distinguish between: $X\beta$ and $X\tilde{\beta}$.

Another consequence of perfect **MC** is that we can not estimate the regression coefficients.

Illustration:

Consider the model in deviations form ($K = 3$):

$$y_i = \beta_2 X_{2i} + \beta_3 X_{3i} + \mathcal{E}_i$$

In relation (2.55), suppose $\alpha_2 = 1$, $\alpha_3 = -5$ and $\alpha_j = 0$ for all other j , i.e., $X_2 = 5X_3$.

$$\alpha_2 X_2 + \alpha_3 X_3 + \dots + \alpha_k X_k = 0 \tag{2.57}$$

We have seen earlier that the **OLS** estimators of β_2 is:

$$\hat{\beta}_2 = \frac{[\sum X_2 y][\sum X_3^2] - [\sum x_3 y][\sum X_2 X_3]}{[\sum X_2^2][\sum x_3^2] - [\sum X_2 X_3]^2}$$

Since we have $X_2 = 5X_3$, we can replace X_2 by $5X_3$:

$$\begin{aligned}\hat{\beta}_2 &= \frac{[\sum (5X_3) y][\sum X_3^2] - [\sum X_3 y][\sum (5X_3) X_3]}{[\sum (5X_3)^2][\sum X_3^2] - [\sum (5X_3) X_3]^2} \\ &= \frac{5[\sum X_3 y][\sum X_3^2] - 5[\sum X_3 y][\sum X_3^2]}{25[\sum X_3^2][\sum X_3^2] - 25[\sum X_3^2]^2} = \frac{0}{0}\end{aligned}$$

Thus, $\hat{\beta}_2$ is indeterminate. It can also be shown that $\hat{\beta}_3$ is indeterminate. Therefore, in the presence of perfect **MC**, the regression coefficients can not be estimated.

2.5.2 Consequences of a high degree of MC (moderate to strong MC)

Consider the case when there is a high degree (moderate to strong) **MC** but not perfect **MC**. What happens to the parameter estimates?

Again consider the model in deviations form ($K = 3$):

$$y_i = \beta_2 X_{2i} + \beta_3 X_{3i} + \mathcal{E}_i$$

There is a high degree of **MC** means that r_{23} , the correlation coefficient between X_2 and X_3 , tends to 1 or -1 (but not equal to ± 1 for this would mean there is perfect **MC**). We can show that the ordinary least squares (**OLS**) estimators of β_2 and β_3 are still **unbiased**, that is,

$$E(\hat{\beta}_j) = \beta_j, j = 2, 3$$

We have seen earlier that the variances of $\hat{\beta}_2$ and $\hat{\beta}_3$ are estimated by:

$$\hat{V}(\hat{\beta}_2) = \frac{\hat{\sigma}^2}{(1 - r_{23}^2)\sum X_{2i}^2} \text{ and } \hat{V}(\hat{\beta}_3) = \frac{\hat{\sigma}^2}{(1 - r_{23}^2)\sum X_{3i}^2}$$

Now, r_{23} tends towards ± 1 :

$\implies r_{23}^2$ approaches to one

$\implies 1 - r_{23}^2$ approaches to zero

\implies both $(1 - r_{23}^2) \sum x_{2i}^2$ and $(1 - r_{23}^2) \sum x_{3i}^2$ approach to zero

\implies both $\hat{V}(\hat{\beta}_2)$ and $\hat{V}(\hat{\beta}_3)$ become very large (or will be inflated)

Particularly, if $r_{23} = \pm 1$, then the variances become infinite.

Recall that to test whether each of the coefficients is significant or not, that is, to test $H_0 : \beta_j = 0$ versus $H_A : \beta_j \neq 0$, the test statistic is:

$$t_j = \frac{\hat{\beta}_j}{\text{s.e.}(\hat{\beta}_j)}, j = 2, 3 \text{ where } \text{s.e.}(\hat{\beta}_j) = \sqrt{\hat{V}(\hat{\beta}_j)}.$$

Thus, under a high degree of **MC**, the standard errors will be inflated and the test statistic will be a very small number. This often leads to incorrectly accepting (not rejecting) the null hypothesis when in fact the parameter is significantly different from zero!

Major implications of a high degree of MC

1. **OLS** coefficient estimates are still unbiased.
2. **OLS** coefficient estimates will have large variances (or the variances will be inflated).
3. There is a high probability of accepting the null hypothesis of zero coefficient (using the t-test) when in fact the coefficient is significantly different from zero.
4. The regression model may do well, that is, R^2 may be quite high.
5. The **OLS** estimates and their standard errors may be quite sensitive to small changes in the data.

Example

Using `data import.csv`, regress imports (Y) on GDP (X_2), stock formation (X_3) and consumption (X_4) for the years 1949 – 1967. check for MC problem?

```
import<-read.csv("import.csv",header = TRUE,sep = ",")
fit<-lm(imports~GDP+stock.formation+consumption,data=import)
summary(fit)
```

##

```
## Call:
## lm(formula = imports ~ GDP + stock.formation + consumption, data = import)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.464 -1.626 -0.113  1.508  3.272
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -19.9824     4.3721   -4.57  0.00037 ***
## GDP              0.0998     0.1936    0.52  0.61393
## stock.formation  0.4467     0.3412    1.31  0.21018
## consumption     0.1488     0.2969    0.50  0.62352
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.39 on 15 degrees of freedom
## Multiple R-squared:  0.975, Adjusted R-squared:  0.97
## F-statistic: 198 on 3 and 15 DF, p-value: 2.78e-12
```

The value of R^2 is close to 1, meaning GDP, stock formation and consumption together explain 0.98% of the variation in imports. Also the F-statistic is significant at the 1% level of significance. Thus, the linear regression model is adequate. However, all of the estimated regression coefficients (save the constant term) are insignificant at the conventional levels of significance. This is an indication that **the standard errors are inflated due to MC**. Since an increase in GDP is often associated with an increase in consumption, they have a tendency to grow up together over time leading to **MC**.

so, let's determine the correlation among them

```
correlation=cor(import)
```

The coefficient of correlation between GDP and consumption is 0.999. Thus, it seems that the

Table 2.8: Correlation Matrix

	year	imports	GDP	stock.formation	consumption
year	1.0000	0.9529	0.9885	0.2452	0.9893
imports	0.9529	1.0000	0.9860	0.3147	0.9859
GDP	0.9885	0.9860	1.0000	0.2681	0.9991
stock.formation	0.2452	0.3147	0.2681	1.0000	0.2660
consumption	0.9893	0.9859	0.9991	0.2660	1.0000

problem of **MC** is due to the joint appearance of these two variables.

2.5.3 Methods of detection of MC

Multicollinearity almost always exists in most applications. So the question is not whether it is present or not; it is a question of degree! Also **MC** is not a statistical problem; it is a data (sample) problem. Therefore, we do not test for "MC"; but measure its degree in any particular sample (using some rules of thumb).

Some of the methods of detecting **MC** are:

1. High R^2 but few (or no) significant t-ratios.
2. High pair-wise correlations among regressor. Note that this is a sufficient but not a necessary condition; that is, small pair-wise correlation for all pairs of regressors does not guarantee the absence of **MC**.

3. Variance inflation factor(VIF)

Consider the regression model:

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_K X_{ki} + \mathcal{E}_i \quad (2.58)$$

The **VIF** of $\hat{\beta}_j$ is defined as:

$$\mathbf{VIF}(\hat{\beta}_j) = \frac{1}{1 - R_j^2}, j = 2, 3, \dots, k$$

where R_j^2 is the coefficient of determination obtained when the X_j variable is regressed on the

remaining explanatory variables (called auxiliary regression). For example, the **VIF** of $\hat{\beta}_2$ is defined as:

$$\mathbf{VIF}(\hat{\beta}_2) = \frac{1}{1 - R_2^2}$$

where R_2^2 is the coefficient of determination of the auxiliary regression:

$$X_{2i} = \alpha_1 + \alpha_3 X_{3i} + \alpha_4 X_{4i} + \dots + \alpha_k X_{ki} + U_i$$

Rule of thumb:

- a) If $\mathbf{VIF}(\hat{\beta}_j)$ exceeds 10, then $\hat{\beta}_j$ is poorly estimated because of **MC** (or the j^{th} regressor variable (x_j) is responsible for **MC**).
- b) (**Klien's rule**) **MC** is trouble some if any of the R_j^2 exceeds the overall R^2 (the coefficient of determination of the regression equation (2.58)).

Example:

Consider the data on imports (Y), **GDP** (X_2), stock formation (X_3) and consumption (X_4) for the years 1949-1967. The coefficient of determination of the auxiliary regression of GDP (X_2) on stock formation (X_3) and consumption (X_4):

$$X_{2i} = \alpha_1 + \alpha_3 X_{3i} + \alpha_4 X_{4i} + U_i$$

```
fit_axui<-lm(GDP~stock.formation+consumption,data=import)
# coefficient of determination
(R_2=summary(fit_axui)$r.squared)
```

```
## [1] 0.9982
```

is $R_2^2 = 0.9982$. The **VIF** of $\hat{\beta}_2$ is thus:

$$VIF(\hat{\beta}_2) = \frac{1}{1 - R_2^2} = \frac{1}{1 - 0.9982} = 556.5817.$$

Since this figure is by far exceeds 10, we can conclude that the coefficient of **GDP** is poorly estimated because of **MC** (or that **GDP** is responsible for **MC**). It can also be shown that $VIF(\hat{\beta}_4) = 556.5817$ indicating that consumption is also responsible for **MC**.

VIF using R package

```
library(car)
vif(fit)
```

```
##                GDP stock.formation      consumption
##                556.58                1.08                555.90
```

Remedial measures

To circumvent the problem of **MC**, some of the possibilities are:

1. Include additional observations maintaining the original model so that a reduction in the correlation among variables is attained.

2. Dropping a variable.

This may result in an incorrect specification of the model (called specification bias). If we consider our example, we expect both **GDP** and Consumption to have an impact on Imports. By dropping one or the other, we have introduced specification bias.

3. A priori information

In the regression model:

$$Y = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + \mathcal{E}$$

suppose X_2 and X_3 are highly collinear. If we have a priori information that $\beta_2 = \lambda \beta_3$, where λ is a

known number, then:

$$\begin{aligned}
 Y &= \beta_1 + \lambda\beta_3X_2 + \beta_3X_3 + \mathcal{E} \\
 &= \beta_1 + \beta_3(\lambda X_2 + X_3) + \mathcal{E} \\
 \text{let, } \lambda X_2 + X_3 &= Z \\
 \text{then, } Y &= \beta_1 + \beta_3Z + \mathcal{E}
 \end{aligned}$$

We can then apply **OLS** to the transformed model and obtain an estimate $\hat{\beta}_3$ of β_3 . An estimator of β_2 is then $\hat{\beta}_2 = \lambda\hat{\beta}_3$.

Drawback: A priori information such as $\beta_2 = \lambda\beta_3$ is rarely available.

4. Use biased (linear) estimators of the β 's . This approach involves trading a little bias for a large reduction in variance. The ordinary ridge regression (**ORR**) is one of such methods.

2.6 Errors in Variables

In applying OLS we assume that the variables (both dependent and independent) are measured without error. The only error is that with the form of disturbance term which represents the influence of relevant explanatory variables that are not included in the model. But sometimes variables are measured subject to error.

Example : Permanent income hypothesis:

$$C_i = \beta_0 + \beta_1 Y_i + U_i \quad (2.59)$$

where Y and C are permanent income and consumption, respectively. But what we actually observe (measure) is Y^* which has a permanent (true) component Y and a transitory component $Y_{\mathcal{E}}$ such that: $Y^* = Y + Y_{\mathcal{E}}$. Similarly, measured consumption expenditure C^* has a permanent (true) component C and a transitory component $C_{\mathcal{E}}$ such that $C^* = C + C_{\mathcal{E}}$. The transitory components represent accidental or chance factors and cyclical variation (such as seasonal effects).

Thus, instead of equation (2.59), what we are estimating is:

$$\begin{aligned}
C_i^* - C_{\mathcal{E}i} &= \beta_0 + \beta_1(Y_i^* - Y_{\mathcal{E}i}) + u_i \\
&\Rightarrow C_i^* = \beta_0 + \beta_1 Y_i^* + (U_i + C_{\mathcal{E}i} - \beta_1 Y_{\mathcal{E}i}) \\
&\Rightarrow C_i^* = \beta_0 + \beta_1 Y_i^* + \eta_i
\end{aligned}$$

where $\eta_i = U_i + C_{\mathcal{E}i} - \beta_1 Y_{\mathcal{E}i}$. This is the case of measurement error in both dependent and independent variables.

2.6.1 Consequences of measurement error in the dependent variable

Suppose our model (in deviation form) is:

$$y_i = \beta x_i + \mathcal{E}_i \quad (2.60)$$

where \mathcal{E}_i represents the random error term. Suppose further that the variable y_i^* instead of y_i is obtained in the measurement process such that:

$$y_i^* = y_i + u_i$$

where u_i is the measurement error. Thus, instead of equation (2.59), the model we are going to estimate is then:

$$y_i^* = y_i + u_i = y_i = \beta x_i + (\mathcal{E}_i + u_i) = y_i = \beta x_i + e_i \quad (2.61)$$

where $e_i = \mathcal{E}_i + u_i$ is the composite error term (model disturbance plus measurement error).

Assumptions :

$$E(\mathcal{E}_i) = E(u_i) = 0$$

$$Var(\mathcal{E}_i) = E(\mathcal{E}_i^2) = \sigma_{\mathcal{E}}^2, Var(u_i) = E(u_i^2) = \sigma_u^2$$

$$cov(u_i, x_i) = cov(\mathcal{E}_i, x_i) = cov(u_i, \mathcal{E}_i) = 0$$

Thus,

$$var(e_i) = var(\mathcal{E}_i + u_i) = var(\mathcal{E}_i) + var(u_i) + 2cov(\mathcal{E}_i, u_i) = \sigma_{\mathcal{E}}^2 + \sigma_u^2$$

The OLS estimator of β from equation (2.61) is:

$$\begin{aligned}\hat{\beta} &= \frac{\sum x_i y_i^*}{\sum x_i^2} = \frac{\sum x_i (\beta x_i + e_i)}{\sum x_i^2} = \beta + \frac{\sum x_i e_i}{\sum x_i^2} \\ &\Rightarrow E(\hat{\beta}) = \beta + \frac{\sum x_i E(e_i)}{\sum x_i^2} = \beta\end{aligned}$$

Thus, $\hat{\beta}$ is unbiased even if there is a measurement error in Y. The variance of $\hat{\beta}$ from equations (2.60) (with no measurement error) is:

$$var(\hat{\beta}) = \frac{\sigma_{\mathcal{E}}^2}{\sum x_i^2}$$

And the variance of $\hat{\beta}$ from equations (2.61) (with measurement error) is:

$$\begin{aligned}var(\hat{\beta}) &= \frac{\sigma_e^2}{\sum x_i^2} \\ &= \frac{\sigma_{\mathcal{E}}^2 + \sigma_u^2}{\sum x_i^2} \\ &= \frac{\sigma_{\mathcal{E}}^2}{\sum x_i^2} + \frac{\sigma_u^2}{\sum x_i^2}\end{aligned}$$

Thus, when there is measurement error in Y, the variance of $\hat{\beta}$ is larger (meaning that the OLS estimator $\hat{\beta}$ of β is no more efficient).

2.7 2.7 Lagged variables

2.8 2.8 Further Model Peculiarities

2.8.1 Stochastic regressors

Suppose the explanatory variable is measured with error, that is, the variable x_i^* instead of x_i is obtained in the measurement process such that:

$$x_i^* = x_i + v_i \implies x_i = x_i^* - v_i$$

Instead of equation (2.60), the model we are going to estimate is then:

$$y_i = \beta(x_i^* - v_i) + \mathcal{E}_i = \beta x_i^* + (\mathcal{E}_i - \beta v_i) = \beta x_i^* + \omega_i \quad (2.62)$$

where $\omega_i = \mathcal{E}_i - \beta v_i$.

Assumptions :

$$E(\mathcal{E}_i) = E(v_i) = 0$$

$$var(\mathcal{E}_i) = E(\mathcal{E}_i^2) = \sigma_{\mathcal{E}}^2, var(v_i) = E(v_i^2) = \sigma_v^2$$

$$cov(v_i, x_i) = cov(\mathcal{E}_i, x_i) = cov(v_i, \mathcal{E}_i) = 0 \Rightarrow E(v_i, x_i) = E(\mathcal{E}_i, x_i) = E(v_i, \mathcal{E}_i) = 0$$

Thus,

$$E(\omega_i) = E(\mathcal{E}_i) - \beta E(v_i) = 0$$

$$E(x_i^*) = x_i + E(v_i) = x_i$$

However, we can not assume that $Cov(\omega_i, x_i^*) = 0$ since:

$$\begin{aligned} Cov(\omega_i, x_i^*) &= E[\omega_i - E(\omega_i)][x_i^* - E(x_i^*)] \\ &= E[\omega_i - 0][x_i^* - x_i] \\ &= E[\mathcal{E}_i - \beta v_i]v_i = E(\mathcal{E}_i v_i) - \beta E(v_i^2) = 0 - \beta \sigma_v^2 = -\beta \sigma_v^2 \neq 0 \end{aligned}$$

Thus, in equation (2.62) the explanatory variable x_i^* is correlated with the error term ω_i . The OLS estimator of β from equation (2.62) is:

$$\hat{\beta} = \frac{\sum x_i^* y_i}{\sum x_i^{*2}} = \frac{\sum x_i^* (\beta x_i^* + \omega_i)}{\sum x_i^{*2}} = \beta + \frac{\sum x_i^* \omega_i}{\sum x_i^{*2}}$$

$$E(\hat{\beta}) = \beta + E \left[\frac{\sum x_i^* \omega_i}{\sum x_i^{*2}} \right] \neq \beta$$

The last inequality holds since $Cov(\omega_i, x_i^*) = E(\omega_i, x_i^*) = -\beta\sigma_v^2 \neq 0$. Thus, $\hat{\beta}$ is a biased estimator of β if there is measurement error in the independent (explanatory) variable X. Expanding the formula for $\hat{\beta}$ above we get:

$$\begin{aligned} \hat{\beta} &= \frac{\sum x_i^* y_i}{\sum x_i^{*2}} = \frac{\sum (x_i + v_i)(\beta x_i + \mathcal{E}_i)}{\sum (x_i + v_i)^2} \\ &= \frac{\beta \sum x_i^2 + \beta \sum x_i v_i + \beta \sum x_i \mathcal{E}_i + \beta \sum \mathcal{E}_i v_i}{\sum x_i^2 + 2 \sum x_i v_i + \sum v_i^2} \end{aligned} \quad (2.63)$$

Note :

1. Let X_n be a sequence of random variables. X_n is said to converge in probability to X_0 if:
 $\lim_{n \rightarrow \infty} P(|X_n - X_0|) = 0$ (written $P \lim X_n = X_0$)
2. If X is a stochastic random variable with zero mean and finite variance, then we have:

$$P \lim \left(\sum_{i=1}^n \frac{X_i^2}{n} \right) = E(X_i^2) = var(X) = \sigma_x^2$$

3. An estimator $\hat{\beta}$ of β is said to be consistent if $p \lim \hat{\beta} = \beta$. That is, by increasing the sample size n, the estimator $\hat{\beta}$ can be made to lie arbitrarily close to the true value β with probability close to one.

From equation (2.63) we have:

$$\begin{aligned}
p \lim \hat{\beta} &= \frac{\beta p \lim(\sum x_i^2/n) + \beta p \lim(\sum x_i v_i/n) + \beta p \lim(\sum x_i \mathcal{E}_i/n) + \beta p \lim(\sum \mathcal{E}_i v_i/n)}{p \lim(\sum x_i^2/n) + p \lim(\sum v_i^2/n) + 2p \lim(\sum x_i v_i/n)} \\
&= \frac{\beta \sigma_x^2 + \beta E(x_i v_i) + E(x_i \mathcal{E}_i) + E(\mathcal{E}_i v_i)}{\sigma_x^2 + E(v_i^2) + 2E(x_i v_i)} \\
&= \frac{\beta \sigma_X^2 + \beta(0) + 0 + 0}{\sigma_X^2 + \sigma_V^2 + 2(0)} \\
&= \frac{\beta \sigma_X^2}{\sigma_X^2 + \sigma_V^2} \\
&= \beta \left(\frac{1}{1 + \sigma_V^2/\sigma_X^2} \right) < \beta
\end{aligned}$$

That is, $p \lim \hat{\beta} \neq \beta$. Thus, if there is measurement error in the independent (explanatory) variable X, then $\hat{\beta}$ is inconsistent. The implication of $p \lim \hat{\beta} < \beta$ is that $\hat{\beta}$ will always underestimate β no matter how large the sample size is.

In the computations above, $E(X_i v_i) = E(X_i \mathcal{E}_i) = E(\mathcal{E}_i v_i) = 0$ (by assumption) and the term σ_x^2 is used to designate $\lim_{n \rightarrow \infty} (\sum_{i=1}^n X_i^2/n) = \lim_{n \rightarrow \infty} (\sum_{i=1}^n (X_i - \bar{X})^2/n)$.

Solutions to errors in variables: the instrumental variable (IV) method

One method which can solve the measurement error problem is the instrumental variables (IV) method.

Given the model: $Y_i = \alpha + \beta X_i + \mathcal{E}_i$, the variable Z is said to be an **instrumental variable** for X if:

$$1. p \lim(\sum Z_i \mathcal{E}_i/n) = 0$$

2. $p \lim(\sum Z_i X_i/n)$ is a finite number different from zero, that is, the correlation (covariance) between Z and X is non-zero as the sample size gets large.

Let us consider the case of measurement error in the independent variable only, that is, the variable x_i^* instead of x_i is obtained in the measurement process such that: $X_i^* = X_i + v_i$. The model we are going to estimate is then:

$$y_i = \beta X_i^* + \omega_i, \text{ with } \omega_i = \mathcal{E}_i - \beta v_i \quad (2.64)$$

We have already seen that the OLS estimator of β is not consistent. To estimate β consistently, we search for an instrumental variable, say, Z that satisfy the above two conditions. These two

conditions are satisfied if z_i is asymptotically uncorrelated with both \mathcal{E}_i and z_i but correlated with x_i^* , that is:

$$\begin{aligned} p \lim(\sum Z_i \mathcal{E}_i / n) &= E(Z_i \mathcal{E}_i) = 0 \\ p \lim(\sum Z_i v_i / n) &= E(Z_i v_i) = 0 \\ p \lim(\sum Z_i x_i^* / n) &\neq 0 \end{aligned}$$

(where $Z_i = Z_i - \bar{Z}$). Then it follows that:

$$p \lim(\sum \omega_i Z_i / n) = E(\omega_i Z_i) = E(\mathcal{E}_i - \beta v_i) Z_i = E(\mathcal{E}_i Z_i) - \beta E(v_i Z_i) = 0$$

An IV estimator of the regression slope in equation (2.64) is:

$$\begin{aligned} \hat{\beta}_{IV} &= \frac{\sum y_i Z_i}{\sum X_i^* Z_i} \\ &= \frac{\sum (\beta X_i^* + \omega_i) Z_i}{\sum X_i^* Z_i} = \beta + \frac{\sum \omega_i Z_i}{\sum X_i^* Z_i} \end{aligned}$$

Taking probability limit, we have:

$$P \lim(\hat{\beta}_{IV}) = \beta + P \lim \left(\frac{\sum \omega_i Z_i}{\sum X_i^* Z_i} \right) = \beta + \frac{P \lim \sum \omega_i Z_i}{P \lim \sum X_i^* Z_i} = \beta, \text{ where } \frac{P \lim \sum \omega_i Z_i}{P \lim \sum X_i^* Z_i} = 0$$

Thus, the IV estimator $\hat{\beta}_{IV}$ is a consistent estimator of β .

2.8.2 Model misspecification

So far we have assumed that the true linear regression relationship is always correctly specified. If any of the assumptions are wrong, there is a specification error. Although some specification errors may have minor implications, others may be more serious.

In developing an empirical model, one is likely to commit one or more of the following specification errors.

- Omission of a relevant variable(s)

- Inclusion of an unnecessary variable(s)
- Adopting the wrong functional form
- Errors of measurement
- Incorrect specification of the stochastic error term

The first four types of error discussed above are essentially in the nature of model specification errors in that we have in mind a “true” model but somehow we do not estimate the correct model. In model misspecification errors, we do not know what the true model is to begin with.

In any econometric investigation, choice of the model is one of the first steps. But, what happen if we use a wrong model and is there ways of assessing whether a model is adequate or not will be covered in this sub topic.

The three **essential features of model choice** are

- ✓ choice of functional form,
- ✓ choice of explanatory variables (regressors) to be included in the model, and
- ✓ whether the multiple regression model assumptions hold.

For choice of functional form and regressors, economic principles and logical reasoning play a prominent and vital role.

2.8.2.1 Underfitting a Model (Omitting a Relevant Variable): It is possible that a chosen model may have important variables omitted. Our economic principles may have overlooked a variable, or lack of data may lead us to drop a variable even when it is prescribed by economic theory.

Suppose the true model is:

$$\text{True model : } Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i \quad (2.65)$$

But the estimated model is:

$$\text{Estimated model : } \hat{Y}_i = \hat{\alpha} + \hat{\beta}_1 X_{1i} \quad (2.66)$$

The estimated model omits a relevant variable X_2 and underspecifies the true relationship. In this case

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_{1i} Y_i}{\sum_{i=1}^n x_{1i}^2} \quad (2.67)$$

Substituting the true model for Y we get

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_{1i} (\alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i)}{\sum_{i=1}^n x_{1i}^2} \quad (2.68)$$

Hence, $E(\hat{\beta}_1) = \beta_1 + \beta_2 b_{12}$ since $E(x_1 u) = 0$ with $b_{12} = \frac{\sum_{i=1}^n x_{1i} X_{2i}}{\sum_{i=1}^n x_{1i}^2}$.

- Note that b_{12} is the regression slope estimate obtained by regressing X_2 on X_1 and a constant. Also, the
- Assume that the estimator of equation (2.65) is given by b_1 and b_2 .

$$\begin{aligned} \text{var}(\hat{\beta}_1) &= E(\hat{\beta}_1 - E(\hat{\beta}_1))^2 = E\left(\frac{\sum_{i=1}^n x_{1i} u_i}{\sum_{i=1}^n x_{1i}^2}\right)^2 \quad (\text{since } \text{var}\left(\beta_2 \frac{\sum_{i=1}^n x_{1i} X_{2i}}{\sum_{i=1}^n x_{1i}^2}\right) = 0 \text{ is non-stochastic}) \\ &= \frac{1}{(\sum_{i=1}^n x_{1i}^2)^2} E\left[\sum_{i=1}^n x_{1i}^2 u_i^2 + \sum_{i \neq j}^n x_{1i} x_{1j} u_i u_j\right] \\ &= \frac{1}{(\sum_{i=1}^n x_{1i}^2)^2} \left[\sum_{i=1}^n x_{1i}^2 E(u_i^2) + \sum_{i \neq j}^n x_{1i} x_{1j} E(u_i u_j)\right] \\ &= \frac{1}{(\sum_{i=1}^n x_{1i}^2)^2} \left[\sum_{i=1}^n x_{1i}^2 \sigma^2 + \sum_{i \neq j}^n x_{1i} x_{1j} \times 0\right] \\ &= \frac{\sigma^2}{\sum x_{1i}^2} \end{aligned}$$

which understates the variance of the estimate of β_1 obtained from the true model, i.e.,

$$b_1 = \frac{\sum \hat{v}_{1i} Y_i}{\sum \hat{v}_{1i}^2} \text{ with}$$

$$\text{var}(b_1) = \frac{\sigma^2}{\sum_{i=1}^n \hat{v}_{1i}^2} = \frac{\sigma^2}{\sum_{i=1}^n x_{1i}^2 (1 - R_1^2)} \geq \text{var}(\hat{\beta}_1) \quad (2.69)$$

In summary, underspecification yields **biased estimates** of the regression coefficients and **under-states** the variance of these estimates. That is

(1) If the left-out, or omitted, variable X_2 is **correlated** with the included variable X_1 , that is, r_{12} , the correlation coefficient between the two variables, is nonzero, $\hat{\alpha}$ and $\hat{\beta}_1$ are **biased** as well as **inconsistent**. That is, $E(\hat{\alpha}) \neq \alpha$ and $E(\hat{\beta}_1) \neq \beta_1$, and the bias does not disappear as the **sample size gets larger**.

(2) If X_1 and X_2 are not correlated then $\hat{\alpha}$ is biased but $\hat{\beta}_1$ is **unbiased**.

2.8.2.2 Inclusion of an Irrelevant Variable (Overfilling a Model): Similarly now assume, the true model is a simple regression with one regressor X_1

$$Y_i = \alpha + \beta_1 X_{1i} + u_i \quad (2.70)$$

but for some reason we fit the following model:

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + v_i \quad (2.71)$$

It can be shown that, from equation (2.70),

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \hat{v}_{1i} Y_i}{\sum_{i=1}^n \hat{v}_{1i}^2} \quad (2.72)$$

where \hat{v}_1 is the OLS residuals of X_1 on X_2 , that is:

$$X_{1i} = \alpha_1 + \alpha_2 X_{2i} + v_{1i}$$

Substituting the true model for Y from (2.65) into (2.67) we get:

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n \hat{v}_{1i}(\alpha + \beta_1 X_{1i} + u_i)}{\sum_{i=1}^n \hat{v}_{1i}^2} \\ &= \frac{\alpha \sum_{i=1}^n \hat{v}_{1i}}{\sum_{i=1}^n \hat{v}_{1i}^2} + \frac{\beta_1 \sum_{i=1}^n \hat{v}_{1i} X_{1i}}{\sum_{i=1}^n \hat{v}_{1i}^2} + \frac{\sum_{i=1}^n \hat{v}_{1i} u_i}{\sum_{i=1}^n \hat{v}_{1i}^2} \\ &= \frac{\beta_1 \sum_{i=1}^n \hat{v}_{1i} X_{1i}}{\sum_{i=1}^n \hat{v}_{1i}^2} + \frac{\sum_{i=1}^n \hat{v}_{1i} u_i}{\sum_{i=1}^n \hat{v}_{1i}^2} \quad (\text{since } \sum_{i=1}^n v_i = 0,)\end{aligned}$$

Again, $X_{1i} = \hat{X}_{1i} + \hat{v}_{1i}$ and $\sum_{i=1}^n \hat{X}_{1i} \hat{v}_{1i} = 0$ implying that

$$\sum_{i=1}^n \hat{v}_{1i} X_{1i} = \sum_{i=1}^n \hat{v}_{1i} (\hat{X}_{1i} + \hat{v}_{1i}) = \sum_{i=1}^n \hat{v}_{1i}^2$$

. Hence,

$$\hat{\beta}_1 = \beta_1 + \frac{\sum_{i=1}^n \hat{v}_{1i} u_i}{\sum_{i=1}^n \hat{v}_{1i}^2} \quad (2.73)$$

and $E(\hat{\beta}_1) = \beta_1$ since \hat{v}_1 is a linear combination of the X's, and $E(X_k u) = 0$ for $k = 1, 2$. Also,

$$\text{var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n \hat{v}_{1i}^2} = \frac{\sigma^2}{\sum_{i=1}^n \hat{x}_{1i}^2 (1 - R_1^2)} \quad (2.74)$$

where $x_{1i} = X_{1i} - \bar{X}_1$ and R_1^2 is the R^2 of the regression of X_1 on X_2 . Using the true model to estimate β_1 , one would get $b_1 = \frac{\sum_{i=1}^n x_{1i} y_i}{\sum_{i=1}^n x_{1i}^2}$ with $E(b_1) = \beta_1$ and $\text{var}(b_1) = \frac{\sigma^2}{\sum_{i=1}^n x_{1i}^2}$. Hence, $\text{var}(\hat{\beta}_1) \geq \text{var}(b_1)$.

Note also that in the **overspecified** model, the estimate for β_2 which has a true value of zero is given by

$$\hat{\beta}_2 = \frac{\sum \hat{v}_{2i} Y_i}{\sum \hat{v}_{2i}^2} \quad (2.75)$$

where \hat{v}_2 is the OLS residual of X_2 on X_1 . Substituting the true model for Y we get

$$\hat{\beta}_2 = \frac{\sum \hat{v}_{2i} u_i}{\sum \hat{v}_{2i}^2} \quad (2.76)$$

since $\sum \hat{v}_{2i}X_{1i} = 0$ and $\sum_{i=1}^n \hat{v}_{2i} = 0$. Hence, $E(\beta_2) = 0$ since \hat{v}_2 is a linear combination of the X 's and $E(X_k u) = 0$ for $k=1,2$. In summary, overspecification still yields unbiased estimates of β_1 and β_2 , but the price is a higher variance.

Example

use data of `educ_inc.csv`. To introduce the **omitted-variable problem**, we consider a sample of married couples such that both husbands and wives work.

The dependent variable is annual family income `FAMINC` defined as the combined income of husband and wife. We are interested in the impact of level of education—both the husband's years of education (`HEDU`) and the wife's years of education (`WEDU`)—on family income.)But for further information, the avriables are `k16`= Number of children less than 6 years old in household, `x5`= `x5` is an artificially generated variable, `x6`= `x6` is an artificially generated variable)

```
rm(list=ls())
educ_inc<-read.csv("educ_inc.csv",header = TRUE,sep = ",")
fit1<-lm(faminc~hedu+wedu,data=educ_inc )
summary(fit1)$coeff
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -5534     11229.5  -0.4928 6.224e-01
## hedu             3132       802.9   3.9002 1.117e-04
## wedu             4523       1066.3   4.2413 2.729e-05
```

The estimated relationship is

$$\widehat{FANINC} = -5534 + 3132 HEDU + 4523 WEDU \quad (2.77)$$

$$(se) = (1.123 \times 10^4) \quad (803) \quad (1066)$$

$$(p - value) = (0.6224) \quad (10^{-4}) \quad (0)$$

We estimate that an additional year of education for the husband will increase annual income by \$3,132, and an additional year of education for the wife will increase income by \$4,523.

What happens if we now incorrectly omit wife's education from the equation?

Table 2.9: Correlation matrix

	faminc	hedu	wedu	k16	x5	x6
faminc	1.0000	0.3547	0.3623	-0.0720	0.2898	0.3514
hedu	0.3547	1.0000	0.5943	0.1049	0.8362	0.8206
wedu	0.3623	0.5943	1.0000	0.1293	0.5178	0.7993
k16	-0.0720	0.1049	0.1293	1.0000	0.1487	0.1595
x5	0.2898	0.8362	0.5178	0.1487	1.0000	0.9002
x6	0.3514	0.8206	0.7993	0.1595	0.9002	1.0000

```
fit2<-lm(faminc~hedu,data=educ_inc )
summary(fit2)$coeff
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    26191      8541.1    3.066 2.304e-03
## hedu           5155       658.5    7.830 3.921e-14
```

The estimated equation becomes

$$FAN\hat{INC} = 2.6191 \times 10^4 + 5155 HEDU \quad (2.78)$$

$$(se) = (8541) \quad (658)$$

$$(p - value) = (0.0023) \quad (0)$$

Relative to (2.77), omitting WEDU leads us to **overstate the effect** of an extra year of education for the husband by about \$2,000. This change in the magnitude of a coefficient is typical of the effect of incorrectly omitting a relevant variable. Omission of a relevant variable (defined as one whose coefficient is nonzero) leads to an estimator that is biased. Naturally enough, this bias is known as omitted-variable bias.

where the correlation between the omitted variable `wedu` and `hedu` is 0.5943.

There are, of course, other variables that could be included in (2.77) as explanators of family income. In the following equation we include `KL6`, the number of children less than six years old. The larger the number of young children, the fewer the number of hours likely to be worked; hence, a lower family income would be expected.

```
fit3<-lm(faminc~hedu+wedu+k16,data=educ_inc )
summary(fit3)$coeff
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -7755     11162.9  -0.6947 4.876e-01
## hedu           3212       796.7   4.0310 6.584e-05
## wedu           4777     1061.2   4.5016 8.727e-06
## k16           -14311     5003.9  -2.8599 4.447e-03
```

$$FAN\hat{I}NC = -7755 + 3212 HEDU + 4777 WEDU + -1.4311 \times 10^4 H16 \quad (2.79)$$

$$(se) = (1.1163 \times 10^4) \quad (797) \quad (1061) \quad (5004)$$

$$(p-value) = (0.4876) \quad (10^{-4}) \quad (0) \quad (0.0044)$$

We estimate that a child under six reduces family income by \$14,311. Notice that compared to (2.77), the coefficient estimates for HEDU and WEDU have not changed a great deal. This outcome occurs because KL6 is not highly correlated with the education variables. From a general modeling perspective, it means that useful results can still be obtained when a relevant variable is omitted if that variable is uncorrelated with the included variables and our interest is on the coefficients of the included variables. (Such instances can arise, for example, if data are not available for the relevant omitted variable.)

Example: Irrelevant variables

To see the effect of irrelevant variables, we add two artificially generated variables X5 and X6 to (2.79). These variables were constructed so that they are correlated with HEDU and WEDU (see Table 2.9) but are not expected to influence family income.

```
fit4<-lm(faminc~hedu+wedu+k16+x5+x6,data=educ_inc )
summary(fit4)$coeff
```

```
##              Estimate Std. Error t value Pr(>|t|)
```

## (Intercept)	-7558.6	11195	-0.6752	0.499948
## hedu	3339.8	1250	2.6717	0.007838
## wedu	5868.7	2278	2.5762	0.010329
## k16	-14200.2	5044	-2.8154	0.005100
## x5	888.8	2242	0.3964	0.692036
## x6	-1067.2	1982	-0.5385	0.590499

The resulting estimated equation is

$$\begin{aligned}
 \hat{FANINC} &= -7559 + 3340 HEDU + 5869 WEDU + -1.42 \times 10^4 H16 + 889 X5 + -1067 X6 \\
 (se) &= (1.1195 \times 10^4) \quad (1250) \quad (2278) \quad (5044) \quad (2242) \quad (1982) \\
 (p-value) &= (0.4999) \quad (0.0078) \quad (0.0103) \quad (0.0051) \quad (0.692) \quad (0.5905)
 \end{aligned}$$

What can we observe from these estimates? First, as expected, the coefficients of X5 and X6 have p-values greater than 0.05. They do indeed appear to be irrelevant variables. Also, the standard errors of the coefficients estimated for all other variables have increased, with p-values increasing correspondingly. The inclusion of irrelevant variables has reduced the precision of the estimated coefficients for other variables in the equation. This result follows because by the Gauss–Markov theorem, the least squares estimator of the correct model is the minimum variance linear unbiased estimator.

2.8.2.3 Model selection criteria: A common feature of the criteria we describe is that they are suitable only for comparing models with the same dependent variable, not models with different dependent variables like y and $\ln(y)$.

These includes:

- ✓ The Adjusted Coefficient of Determination (\bar{R}^2), $\bar{R}^2 = 1 - \frac{ESS/(n-k)}{TSS/(n-1)}$
- ✓ The Akaike information criterion (AIC), $AIC = \ln\left(\frac{ESS}{n}\right) + \frac{2k}{n}$
- ✓ The Schwarz criterion (SC), also known as the Bayesian information criterion (BIC),

$$SC = \ln\left(\frac{ESS}{n}\right) + \frac{k \times \ln(n)}{n}$$

Using these last two criteria, the model with the smallest AIC, or the smallest SC, is preferred.

To get values of the more general versions of these criteria based on maximized values of the likelihood function, you need to add $1 + \ln(2\pi)$ to AIC and SC equation.

Example: model selection

```
# the 4 fits were
fit1<-lm(faminc~hedu,data=educ_inc )
fit2<-lm(faminc~hedu+wedu,data=educ_inc )
fit3<-lm(faminc~hedu+wedu+k16,data=educ_inc )
fit4<-lm(faminc~hedu+wedu+k16+x5+x6,data=educ_inc )

##### AIC,SC and R2 calculation
#AIC from Fit 1
r1=resid(fit1); ess1=sum(r1^2);log_value1=log(ess1/nrow(educ_inc))
(aic_1=log_value1+(2*2/nrow(educ_inc)))

## [1] 21.26

#AIC from fit2
r2=resid(fit2); ess2=sum(r2^2); log_value2=log(ess2/nrow(educ_inc))
aic_2=log_value2+(2*3/nrow(educ_inc))
#AIC from Fit3
r3=resid(fit3); ess3=sum(r3^2);log_value3=log(ess3/nrow(educ_inc))
aic_3=log_value3+(2*4/nrow(educ_inc))
# AIC from Fit4
r4=resid(fit4); ess4=sum(r4^2);log_value4=log(ess4/nrow(educ_inc))
aic_4=log_value4+(2*6/nrow(educ_inc))

#####
```



```

# SC from Fit 1

r1=resid(fit1); ess1=sum(r1^2); log_value1=log(ess1/nrow(educ_inc))
(sc_1=log_value1+(2*log(nrow(educ_inc)))/nrow(educ_inc))

## [1] 21.28

# SC from Fit 2
r2=resid(fit2); ess2=sum(r2^2); log_value2=log(ess2/nrow(educ_inc))
sc_2=log_value2+(3*log(nrow(educ_inc)))/nrow(educ_inc))

# SC from Fit 3
r3=resid(fit3); ess3=sum(r3^2); log_value3=log(ess3/nrow(educ_inc))
sc_3=log_value3+(4*log(nrow(educ_inc)))/nrow(educ_inc))

# SC from Fit 4
r4=resid(fit4); ess4=sum(r4^2); log_value4=log(ess4/nrow(educ_inc))
sc_4=log_value4+(6*log(nrow(educ_inc)))/nrow(educ_inc))

#
options(digits = 6)

# To combine into table form, so as make vector
scc=c(sc_1,sc_2,sc_3 ,sc_4)
aic=c(aic_1,aic_2,aic_3,aic_4)
R_square=c(summary(fit1)$r.squared,summary(fit2)$r.squared,
            summary(fit3)$r.squared,summary(fit4)$r.squared)
R_sq_adj=c(summary(fit1)$adj.r.squared,summary(fit2)$adj.r.squared,
            summary(fit3)$adj.r.squared,summary(fit4)$adj.r.squared)
explanatory_variables=c("HEDU", "HEDU;WEDU", "HEDU;WEDU;KL6",
                        "HEDU;WEDU;KL6;X5;X6")

```

From Table 2.10:

- 1) Adding more variables always increases the R^2 , whether they are relevant or not.
- 2) The \bar{R}^2 increases when relevant variables are added, but declines in the last case when the

Table 2.10: Goodness-of-Fit and Information Criteria for Family Income Example

variables	R_square	R_sq_adj	AIC	Sc
HEDU	0.1258	0.1237	21.262	21.281
HEDU;WEDU	0.1613	0.1574	21.225	21.253
HEDU;WEDU;KL6	0.1772	0.1714	21.211	21.248
HEDU;WEDU;KL6;X5;X6	0.1778	0.1681	21.219	21.276

irrelevant variables X5 and X6 are added.

- 3) The AIC and SC are smallest for the model with variables HEDU, WEDU, and KL6. Thus, in this case, but not necessarily in general, maximizing \bar{R}^2 , minimizing AIC, and minimizing SC all lead to selection of the same model.

2.8.3 Qualitative variables

Indicator variables are used to account for qualitative factors in econometric models. Also called **dummy, binary or dichotomous** variables, because they take just **two values**, usually **one or zero**, to indicate the **presence or absence of a characteristic** or to indicate whether a condition is true or false.

Generally, we define an indicator variable D as

$$D = \begin{cases} 1, & \text{if characteristic is present} \\ 0, & \text{if characteristic is not present} \end{cases}$$

let's take an example about real estate economics to predict a house price for buyers and sellers. Price of a house can be modeled as a function some characteristics such as its size, location, number of bedrooms, age, and so on. The idea is to break down a good into its component pieces, and then estimate the value of **each characteristic**.

But now let us assume that the size of the house, measured in square feet, SQFT, is the only relevant variable in determining house price, PRICE, as follow.

$$PRICE = \beta_1 + \beta_2 SQFT + \mathcal{E}$$

where β_1 = the value of the land alone, and

β_2 = the value of an additional square foot of living area.

In real estate modeling price, location is factor. So, how can we take into account the effect of a property's being in a desirable neighborhood(location), such as one near a university or others?

Here, location is a **qualitative** characteristic of a house and we account to the house price model by defining an indicator or dummy variable as follow:

$$D = \begin{cases} 1, & \text{if property is in the desirable neighborhood} \\ 0, & \text{if property is not in the desirable neighborhood} \end{cases}$$

Indicator variables can be used to capture changes in the model intercept, or slopes, or both. We consider these possibilities in turn.

2.8.3.1 Intercept indicator variables: To assess whether two regression models are important for modeling the house price in the two locations(near the university and not), we add an indicator variable to modify the regression model intercept parameter as follows.

$$PRICE = \beta_1 + \delta D + \beta_2 SQFT + \mathcal{E}$$

where δ = a location premium, the difference in house price due to the houses being located in the desirable neighborhood.

The effect of the inclusion of an indicator variable D into the regression model is shown as:

$$E(PRICE) = \begin{cases} \beta_1 + \delta \times 1 + \beta_2 SQFT = (\beta_1 + \delta) + \beta_2 SQFT, & \text{if } D=1 \\ \beta_1 + \delta \times 0 + \beta_2 SQFT = \beta_1 + \beta_2 SQFT, & \text{if } D=0 \end{cases}$$

D=1 if house is in the desirable neighborhood

D=0 if house is not in the desirable neighborhood

An indicator variable like D that is incorporated into a regression model to capture a shift in the intercept as the result of some qualitative factor is called an **intercept indicator variable, or an**

intercept dummy variable(as shown in Figure 2.14).

Use a **t-test** to test whether the neighborhood effect on house price is **statistically significant** and hence a model with an intercept indicator variable or not is to use.

$H_0 : \delta = 0$ (No location effect in price of house)

$H_1 : \delta \neq 0$ (there is a location effect in price of house)

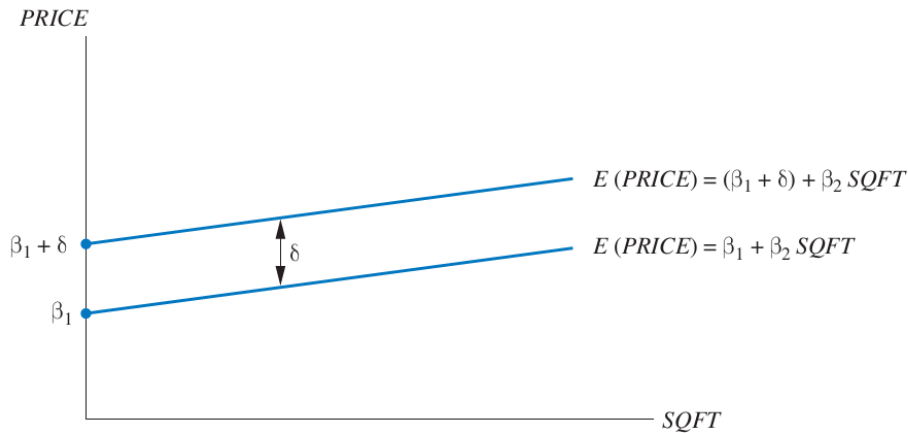


Figure 2.14: An intercept indicator variable.

2.8.3.2 Slope-indicator variables To allow a difference in slope of the relationship, we have to include an **interaction variable** by combining an indicator variable and other continuous variable as follow:

$$PRICE = \beta_1 + \beta_2 SQFT + \gamma(SQFT \times D) + \mathcal{E}$$

where γ = a difference or a change due to the interaction effect of location and size on house price. The **new variable** $(SQFT \times D)$ is the product of house size and the indicator variable, and is called an **interaction variable**, as it captures the interaction effect of location and size on house price. Alternatively, it is called a slope-indicator variable or a **slope dummy** variable, because it

allows for a change in the slope of the relationship.

Examining the regression function for the two different locations best illustrates the effect of the inclusion of the slope-indicator variable into the economic model,

The effect of the inclusion of the slope-indicator variable in the house price model between the two locations is shown:

$$E(PRICE) = \beta_1 + \beta_2 SQFT + \gamma(SQFT \times D) = \begin{cases} \beta_1 + (\beta_2 + \gamma)SQFT, & \text{if } D=1 \\ \beta_1 + \beta_2 SQFT, & D=0 \end{cases}$$

In the desirable neighborhood, the price per additional square foot of a home is $(\beta_2 + \gamma)$ but in the other locations is β_2 , and look (Figure 2.15(a)) for the effect of inclusion of slope-indicator variable.

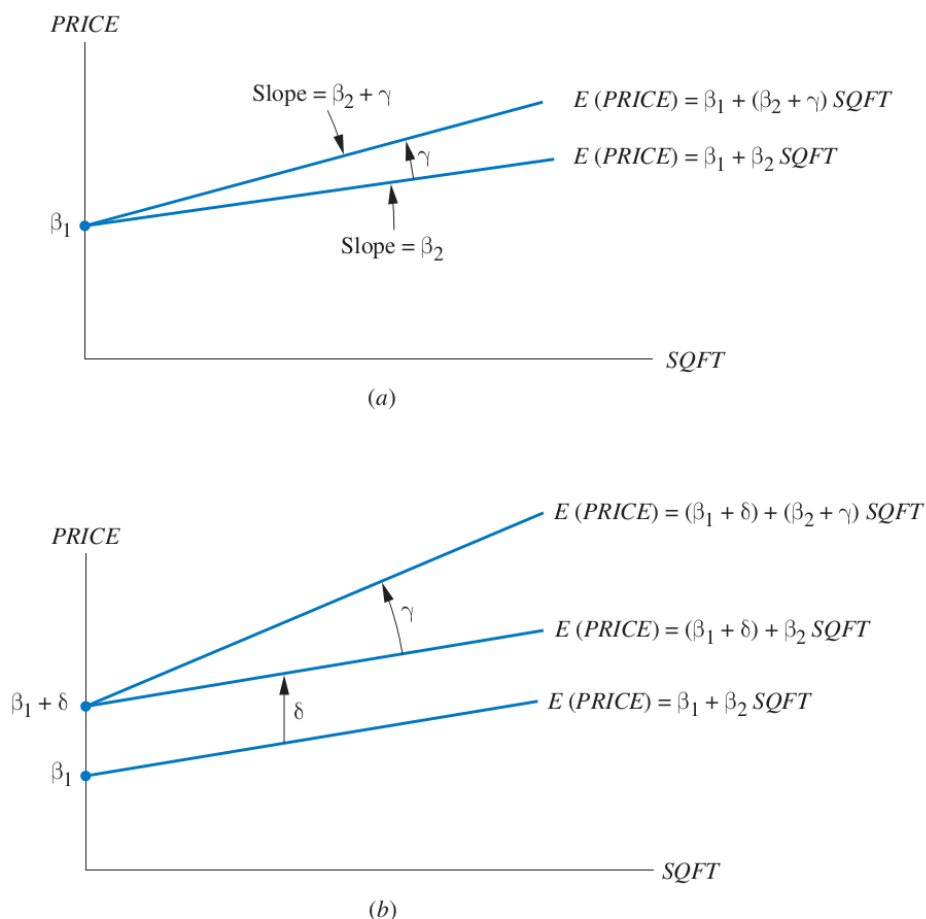


Figure 2.15: (a) A slope-indicator variable. (b) Slope- and intercept-indicator variables.

A test of the hypothesis that the value of an additional square foot of living area is the same in the two locations is carried out by testing the null hypothesis using t-test

$$H_0 : \gamma = 0 \text{ Vs } H_1 : \gamma \neq 0$$

If we assume that house location affects **both the intercept and the slope**, then both effects can be incorporated into **a single model**. The resulting regression model is

$$PRICE = \beta_1 + \delta D + \beta_2 SQFT + \gamma(SQFT \times D) + \varepsilon$$

In this case the regression functions for the house prices in the two locations are

$$E(PRICE) = \beta_1 + \delta D + \beta_2 SQFT + \gamma(SQFT \times D) = \begin{cases} (\beta_1 + \delta) + (\beta_2 + \gamma)SQFT, & \text{if } D=1 \\ \beta_1 + \gamma SQFT, & D=0 \end{cases}$$

In Figure 2.15b we depict the house price relations assuming that $\delta > 0$ and $\gamma > 0$.

Example: The university effect on house prices

A real estate economist collects information on 1000 house price sales from two similar neighborhoods, one called “University Town” bordering a large state university, and one a neighborhood about three miles from the university. A few of the observations are shown in (Table 2.11). The complete data file is **utown.csv**.

House prices are given in \$1,000; size (SQFT) is the number of hundreds of square feet of living area. For example, the first house sold for \$205,452 and has 2346 square feet of living area. Also recorded are the house AGE(in years),

- location (UTOWN=1 for homes near the university, 0 otherwise),
- whether the house has a pool (POOL=1 if a pool is present, 0 otherwise) and
- whether the house has a fireplace (FPLACE=1 if a fireplace is present, 0 otherwise).

Table 2.11: Representative Real Estate Data Values

price	sqft	age	utown	pool	fplace
205.5	23.46	6	0	0	1
185.3	20.03	5	0	0	1
248.4	27.77	6	0	0	0
154.7	20.17	1	0	0	0
221.8	26.45	0	0	0	1
199.1	21.56	6	0	0	1

Table 2.12: House Price Equation Estimates

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	24.5000	6.1917	3.957	0.0001
utown	27.4530	8.4226	3.259	0.0012
sqft	7.6122	0.2452	31.048	0.0000
I(sqft * utown)	1.2994	0.3320	3.913	0.0001
age	-0.1901	0.0512	-3.712	0.0002
pool	4.3772	1.1967	3.658	0.0003
fplace	1.6492	0.9720	1.697	0.0901

The economist specifies the regression equation as

$$PRICE = \beta_1 + \delta_1 UTOWN + \beta_2 SQFT + \gamma(SQFT \times UTOWN) + \beta_3 AGE + \delta_2 POOL + \delta_3 FPLACE + \mathcal{E}$$

Solution

```
house_price<-read.csv("utown.csv",header = TRUE,sep=",")

fit<-lm(price~utown+sqft+I(sqft*utown)+age+pool+fplace,data=house_price)
#summary(fit)
```

The goodness of fit statistic

Note that POOL and FPLACE are **intercept dummy variables**. By introducing these variables we are asking whether, and by how much, these features change house price. Because these variables **stand alone**, and are not interacted with SQFT, we are assuming that they affect the regression intercept, but not the slope.

Table 2.13: The goodness of fit statistic and F

	x
R_square	0.8706
F	1113.1828

The slope-indicator variable is $SQFT \times UTOWN$. Based on the t-test significance all have $p\text{-value} < 0$

In particular, based on these t-tests, we conclude that houses near the university have a significantly higher **base price**, and that their price per additional square foot is significantly higher than in the comparison neighborhood

The estimated regression function for the houses near the university is:

$$\hat{PRICE} = (24.5 + 27.453) + (7.612 + 1.299)SQFT - 0.19AGE + 4.377POOL + 1.649FPLACE$$

$$\hat{PRICE} = 51.953 + 8.911SQFT - 0.19AGE + 4.377POOL + 1.649FPLACE$$

For houses in other areas, the estimated regression function is

$$\hat{PRICE} = 24.5 + 7.612SQFT - 0.19AGE + 4.377POOL + 1.649FPLACE$$

‘ Based on the regression results in Table 10, we estimate that

- ✓ The location premium for lots near the university is \$27,453
- ✓ The change in expected price per additional square foot is \$89.12 for houses near the university and \$76.12 for houses in other areas
- ✓ Houses depreciate \$\$\$190.10 per year
- ✓ A pool increases the value of a home by \$\$\$4,377.20
- ✓ A fireplace increases the value of a home by \$\$\$1,649.20

2.8.3.3 Tests of structural stability: Suppose we are interested in estimating a simple saving function that relates domestic household savings (S) with gross domestic product (Y) for a certain

country. Suppose further that, at a certain point of time, a series of economic reforms have been introduced. The hypothesis here is that such reforms might have considerably influenced the savings-income relationship, that is, the relationship between savings and income might be different in the post-reform period as compared to that in the pre-reform period. If this hypothesis is true, then we say **a structural change** has happened. How do we check if this is so?

1. Chow's test

One approach for testing the presence of structural change (structural instability) is by means of Chow's test. The steps involved in this procedure are as follows:

a) Estimate the regression equation:

$$S_i = \alpha + \beta Y_i + \mathcal{E}_i, \dots i = 1, 2, 3, \dots, n \quad (2.80)$$

for the whole period (pre-reform plus post-reform periods) and find the error sum of squares (ESS_R).

b) Estimate equation (2.80) using the available data in the pre-reform period (say, of size n_1), that is, estimate the model:

$$S_i = \alpha_1 + \beta_1 Y_i + \mathcal{E}_i, \dots i = 1, 2, 3, \dots, n_1$$

and find the error sum of squares (ESS_1).

c) Estimate equation (2.80) using the available data in the post-reform period (say, of size n_2), that is, estimate the model:

$$S_i = \alpha_2 + \beta_2 Y_i + \mathcal{E}_i, \dots i = 1, 2, 3, \dots, n_2$$

and find the error sum of squares (ESS_2).

d) Calculate: $ESS_{UR} = ESS_1 + ESS_2$

e) Calculate the Chow test statistic:

$$F = \frac{\frac{(ESS_R + ESS_{UR})}{K}}{\frac{ESS_{UR}}{(n_1 + n_2 - 2k)}}$$

Table 2.14: Saving data

year	saving	GDP
1980	27136	401128
1981	31355	425073
1982	34368	438080
1983	38587	471742
1984	46063	492077
1985	54167	513990

where K is the number of estimated regression coefficients.

f) *Decisionrule*: Reject the null hypothesis of identical intercepts and slopes for the pre-reform and post-reform periods, that is,

$$H_0 : \left\{ \begin{matrix} \alpha_1 = \alpha_2 \\ \beta_1 = \beta_2 \end{matrix} \right\}$$

$$if : F > F_\alpha(k, n_1 + n_2 - 2k)$$

where $F_\alpha(k, n_1 + n_2 - 2k)$ is the critical value from the F-distribution with K (in our case $K = 2$) and $n_1 + n_2 - 2k$ degrees of freedom for a given significance level α .

Note that rejecting H_0 means there is a structural change.

Example: The following data is on domestic household savings (S) and gross domestic product (Y) for India for the period 1980 to 2002 (see Table 2.14).

```
total<-read.csv("saving.csv",header = TRUE,sep = ",")
```

Let's separate the data to pre-reform and post-reform

```
pre_reform<-total[1:12,]
post_reform<-total[13:23,]
```

Now let's regress saving on GDP for pre reform and post reform as follow(step (b),(c))

```
fit1=lm(saving~GDP,data=pre_reform)
fit2=lm(saving~GDP,data=post_reform)
resd1=resid(fit1)
resd2=resid(fit2)
```

```
## Calculate sum square of error of each model
```

```
options(digits = 14)
( ESS1=sum(resd1^2) )
```

```
## [1] 644994361.86529
```

```
( ESS2=sum(resd2^2) )
```

```
## [1] 2736652790.4339
```

Step a

```
#lm of the total, assuming that no difference(Pooled)
```

```
fit_tot=lm(saving~GDP,data=total)
resd_tot=resid(fit_tot)
( ESS_R=sum(resd_tot^2) )
```

```
## [1] 13937337067.461
```

Step d and e

```
(ESS_UR=ESS1+ESS2)
```

```
## [1] 3381647152.2992
```

```
k=2
```

```
n1=nrow(pre_reform)
n2=nrow(post_reform)
```

```
# F-test (( ) / ( ) ) (( ) / ( ) )
```

```
(F_cal= ( (ESS_R-ESS_UR) / k ) / ( (ESS_UR) / (n1+n2-2*k) ) )
```

```
## [1] 29.653908192597
```

```
# F_tab
```

```
(qf(p=0.05,df1 = k,df2=n1+n2-2*k,lower.tail = F) )
```

```
## [1] 3.5218932605788
```

```
# p-value
```

```
(pf(F_cal, df1=n1-2, df2=n1+n2-2, lower.tail = F) )
```

```
## [1] 3.7619822496523e-10
```

Decision: Since the calculated value of F exceeds the tabulated value, we reject the null hypothesis of identical intercepts and slopes for the pre-reform and post-reform periods at the 5% level of significance. Thus, we can conclude that there is a structural change.

Drawback: Chow's test does not tell us whether the difference (change) is in the slope only, in the intercept only or in both the intercept and the slope.

Using dummy variables Write the savings function as:

$$S_t = \beta_0 + \beta_1 D_t + \beta_2 Y_t + \beta_3 (D_t Y_t) + U_t \quad (2.81)$$

where S_t is household saving at time t, Y_t is GDP at time t and:

$$D_t = \begin{cases} 0 \rightarrow \text{pre-reform period} \\ 1 \rightarrow \text{post-reform period} \end{cases}$$

Here β_3 is the differential slope coefficient indicating how much the slope coefficient of the pre-reform period savings function differs from the slope coefficient of the savings function in the post reform period. Observe that:

$$E\left(\frac{S_t}{D_t} = 0, Y_t\right) = \beta_0 + \beta_2 Y_t$$
$$E\left(\frac{S_t}{D_t} = 1, Y_t\right) = (\beta_0 + \beta_1) + (\beta_2 + \beta_3) Y_t$$

If β_1 and β_3 are both statistically significant as judged by the t-test, then the pre-reform and post-reform regressions differ in both the intercept and the slope. However, if only β_1 is statistically significant, then the pre-reform and post-reform regressions differ only in the intercept (meaning the marginal propensity to save (MPS) is the same for pre-reform and post-reform periods). Similarly, if

Table 2.15: Dummy variable regression result

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-1.337e+05	2.126e+04	-6.289	0
D	-2.286e+05	3.078e+04	-7.429	0
GDP	3.750e-01	3.860e-02	9.717	0
D:GDP	3.385e-01	4.410e-02	7.673	0

only β_3 is statistically significant, then the two regressions differ only in the slope (MPS).

Estimating equation (2.81) using the above data by OLS yields the following result.

```
options(digits = 4)
# Creating indicator variable
total$D[total$year<=1991]=0 # pre=reform is taken as base
total$D[total$year>1991]=1

fit=lm(saving~D+GDP+D:GDP,data = total) # lm(saving~D*GDP,data = total)
#model.matrix(fit)
```

From Table 2.15, we can see that both the differential intercept coefficient $\hat{\beta}_1 = -2.2863 \times 10^5$ and differential slope coefficient $\hat{\beta}_3 = 0.3385$ are statistically significant. Thus, the savings-income relationship for the two periods is different.

```
## difference of intercept
summary(fit2)$coeff[1]-summary(fit1)$coeff[1] # summary(fit)$coeff[2]

## [1] -228629

# difference in slope
summary(fit2)$coeff[2]-summary(fit1)$coeff[2] #summary(fit)$coeff[4]

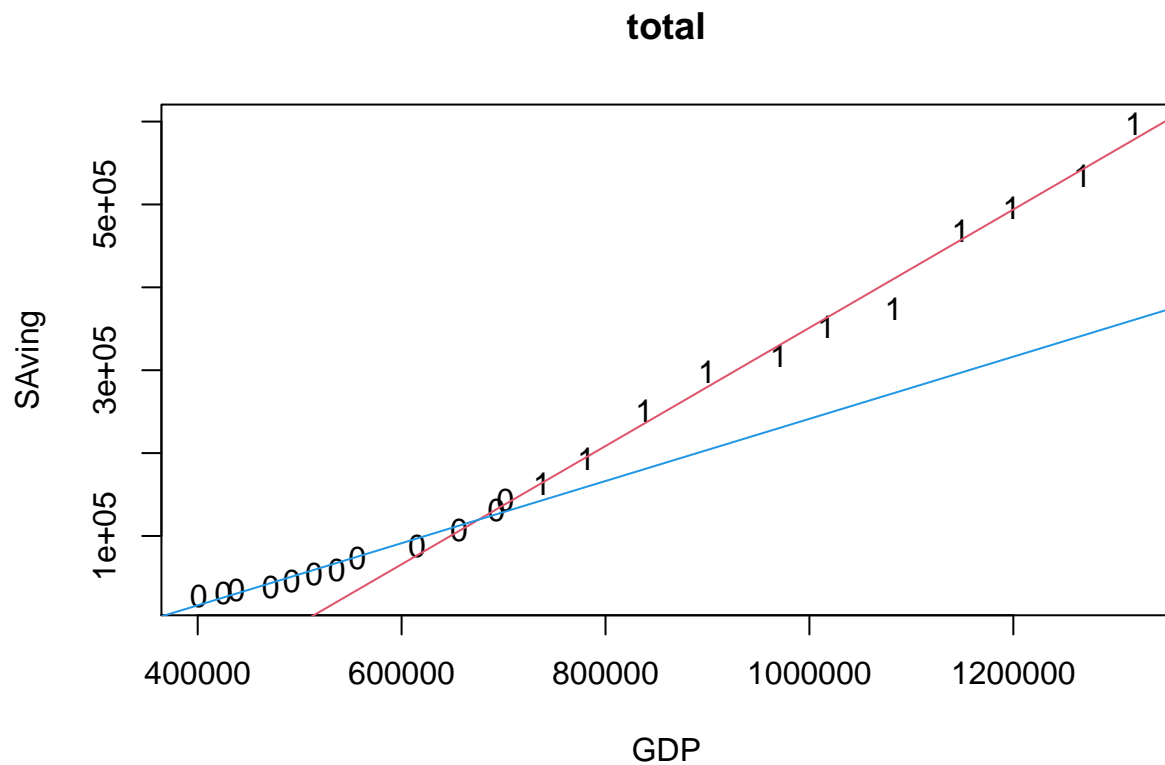
## [1] 0.3385

lapply(split(total,total$D),summary)
```

```

plot(total$GDP,total$saving,type="n",xlab="GDP",ylab="SAving",
      main="total")
text(total$GDP,total$saving,total$D)
abline(summary(fit)$coef[1]+summary(fit)$coef[2],summary(fit)$coef[3]+
        summary(fit)$coef[4],col=2)
abline(summary(fit)$coef[1],summary(fit)$coef[3],col=4) # base is pre

```



The graph also clearly shows difference both in intercept and slope.

3 Non-linear Models

3.1 Introduction

The major emphasis of the course is on linear regression models, that is, models that are linear in the parameters and/or models that can be transformed so that they are linear in the parameters. But in this chapter we will see models that are nonlinear in the parameters.

We say a linear regression model if it is linear in the model parameters but may or may not linear in the variables. A regression model is called nonlinear, if the **derivatives** of the model with respect to the model parameters depends on one or more parameters. Nonlinear regression model is a nonlinear in the parameters and may or may not linear in the variables.

Assume we have a model:

$$Y = e^{\alpha + \beta X + \mathcal{E}} \quad (3.82)$$

Now, let's apply on both sides a natural logarithm.

$$\ln Y = \ln e^{\alpha + \beta X + \mathcal{E}}$$

$$\ln Y = (\alpha + \beta X + \mathcal{E}) \ln e = \alpha + \beta X + \mathcal{E}$$

Let $Y^* = \ln Y$, then we can re-write the above equation as $Y^* = \alpha + \beta X + \mathcal{E}$. Here Y^* is linear in the the model parameters, α and β . Therefore, there are models which may look nonlinear in the parameters but are **inherently** or **intrinsically** linear because with suitable transformation they can be made linear in the parameter regression models.

3.2 Intrinsically Linear Models

A non-linear model with respect to the variables but linear with respect to the parameters to be estimated. Suitable transformations of data can frequently (not always) be found that will reduce a theoretical nonlinear model to a linear form. Linearizing may require transforming both the independent and dependent variable. Hence, model in 3.1 is an example of intrinsically linear model. The basic common characteristic of such models is that they can be converted into ordinary

linear models by suitable transformation of variables – and such transformation amounts to nothing more than re-labeling one or more of the variables. While it may be useful, at times, to transform a model of this type so that it can be easily fitted.

consider the famous Cobb–Douglas (C–D) production function.

$$Y_i = \beta_1 X_{2i}^{\beta_2} X_{3i}^{\beta_3} \mathcal{E}^{u_i} \quad (3.83)$$

Letting Y =output, X_2 =labor input, and X_3 =capital input, we will write this function in three different ways

3.3 Intrinsically Non-linear Models

Not all functions are linearizable, nor in some cases it is desirable to transform to linearity.

Let $Y = \frac{\alpha}{\alpha - \beta} [e^{-\beta X} - e^{\alpha X}] + \mathcal{E}$ is impossible to convert into a form of linear in parameters, and we will call **intrinsically nonlinear model**. From now on when we talk about a nonlinear regression model, we mean that it is intrinsically nonlinear.

Exercise

1) Are the following models linear regression models? Why or why not?

a)

$$Y_i = e^{\beta_1 + \beta_2 X_i + \mathcal{E}}$$

b)

$$Y_i = \frac{1}{1 + e^{\beta_1 + \beta_2 X_i + \mathcal{E}}}$$

c)

$$\ln Y_i = \beta_1 + \beta_2 (1/X_i) + \mathcal{E}$$

Qualitative choice analysis

Qualitative choice models may be used when a decision maker faces a choice among a set of alternatives meeting the following criteria:

. The number of choices if finite

. The choices are mutually exclusive (the person chooses only one of the alternatives) . The choices are exhaustive (all possible alternatives are included)

The first criterion is a binding one. We can always refine the available choices so that they can satisfy the last two criteria. Throughout our discussion we shall restrict ourselves to cases of qualitative choice where the set of alternatives is binary. For the sake of convenience the dependent variable is given a value of 0 or 1.

Example : Suppose the choice is whether to work or not. The discrete dependent variable we are working with will assume only two values 0 and 1:

$$Y_i = \begin{cases} 1 \rightarrow \text{if } i^{\text{th}} \text{ individual is working/seeking work} \\ 0 \rightarrow \text{if } i^{\text{th}} \text{ individual is not working} \end{cases}$$

where $i = 1, 2, \dots, n$. The independent variables (called factors) that are expected to affect an individual's choice may be X_1 = age, X_2 = marital status, X_3 = gender, X_4 =education, etc. These are represented by a matrix X. If we have k factors, the vector of factors for the i^{th} individual is given by:

$$X_i(x_{1i}, x_{2i}, x_{3i}, \dots, x_{ki}) \quad i = 1, 2, \dots, n$$

The regression approach

The economic interpretation of discrete choice models is typically based on the principle of utility maximization leading to the choice of, say, A over B if the utility of A exceeds that of B. For example, let U^1 be the utility from working/seeking work and let U^0 be the utility from not working. Then an individual will choose to be part of the labour force if $U^1 - U^0 > 0$, and this decision depends on a number of factors X.

The probability that the i^{th} individual chooses alternative 1 (i.e. works/seeks work) given his/her individual characteristics X_i is:

$$P_i = \text{Prob}(Y_i = 1 | X_i) = \text{Prob}[(U^1 - U^0)_i > 0] = G(X_i, \beta)$$

The vector of parameters $\beta (\beta = (\beta_1, \beta_2, \dots, \beta_k)')$ measures the impact of changes in X (say, age,

marital status, gender, education, etc) on the probability of labour force participation.

And the probability that the i^{th} individual chooses alternative 0 (i.e. not to work) is given by:

$$Prob(Y_i = 0|X_i) = 1 - P_i = 1 - Prob[(U^1 - U^0)_i > 0] = 1 - G(X_i, \beta)$$

Here P_i is called the **response probability** and $(1 - P_i)$ is called the **non-response probability**

The mean response of the i^{th} individual given his/her individual characteristics X_i is:

$$E(Y_i|X_i) = 1 \times Prob(Y_i = 1|X_i) + 0 \times Prob(Y_i = 0|X_i) = 1 \times G(X_i, \beta) + 0 \times G(X_i, \beta) = G(X_i, \beta)$$

The problem is thus to choose the appropriate form of $G(X_i, \beta)$.

Case 1: The linear probability model

The linear probability model defines $G(X_i, \beta)$ as: $G(X_i, \beta) = X_i\beta$

The regression model is thus:

$$Y_i = X_i\beta + \mathcal{E}_i \quad (3.84)$$

This is the usual linear regression model. The drawbacks of this model are:

1. The right hand side of equation (3.84) is a combination of discrete and continuous variables while the left hand side variable is discrete.
2. Usually we arbitrarily (or for convenience) use 0 and 1 for Y_i . If we use other values for Y_i , say 3 and 4, β will also change even if the vector of factors X_i remains unchanged.
3. \mathcal{E}_i assumes only two values:

if $Y_i = 1$ then $\mathcal{E}_i = 1 - X_i\beta$ (with prob. P_i) if $Y_i = 0$ then $\mathcal{E}_i = -X_i\beta$ (with prob. $1 - P_i$)

Consequently, \mathcal{E}_i is not normally distributed but rather has a discrete (binary) probability distribution defined as:

\mathcal{E}_i	Probability
$1 - X_i\beta$	P_i
$-x_i\beta$	$1 - P_i$

4. The expectation (mean) of \mathcal{E}_i conditional on the exogenous variables X_i is (from the above

table):

$$E(\mathcal{E}_i|x_i) = (1 - X_i\beta)P_i + (-x_i\beta)(1 - P_i) = P_i - X_i\beta$$

Setting this mean to zero as in the classical regression analysis means:

$$E(\mathcal{E}_i|x_i) = 0 \Rightarrow P_i = X_i\beta$$

So the original model (1) becomes:

$$Y_i = P_i + \mathcal{E}_i \Rightarrow \mathcal{E}_i = Y_i - P_i$$

That is, the binary (discrete) disturbance term \mathcal{E}_i is equal to the difference between a binary variable Y_i and a continuous response probability P_i . Clearly this does not make sense.

5. We know that the probability of an event is always a number between 0 and 1 (inclusive). But here we can see that:

$$P_i = Prob(Y_i = 1|X_i) = X_i\beta$$

i.e., P_i can take on any value (even negative numbers) leading to nonsense probabilities.

The latent regression approach

The outcome (or observed occurrence) of a discrete choice may be considered to be an indicator of an underlying, unobservable continuous variable which may be called ‘propensity to choose a given alternative’. Such a variable is characterized by the existence of a threshold, where crossing a threshold means switching from one alternative to another. For instance, a married woman’s propensity to join the labour force may be directly related to the wage that she may receive in the market, which in turn may depend on her level of education, experience, etc. Whether she actually joins the work force or not is likely to depend on whether her market wage does or does not exceed her threshold or ‘reservation’ wage.

Example: Let an individual’s propensity to enter the labour force (or to work) be an unobservable

(latent) variable Y_i^* such that:

$$Y_i^* = X_i\beta + \mathcal{E}_i = X_{1i}\beta_1 + X_{2i}\beta_2 + \dots + X_{ki}\beta_k + \mathcal{E}_i$$

What we actually observe here is the individual's decision to work or not (Y_i) and the set of individual factors such as age, sex, marital status, level of education, experience, etc ($X_i(X_{1i}, X_{2i}, X_{3i}, \dots, X_{ki})$) where:

$$Y_i = \begin{cases} 1 \rightarrow & \text{if } Y_i^* > c \\ 0 \rightarrow & \text{if } Y_i^* \leq c \end{cases}$$

where c is a floor for the propensity variable (or a threshold). The probability of labour force participation of the i^{th} individual given his/her personal background X_i is:

$$Prob(Y_i = 1|X_i) = Prob(Y_i^* > c|X_i) = Prob(X_i\beta + \mathcal{E}_i > c) = Prob(\mathcal{E}_i > c - X_i\beta|X_i)$$

Similarly,

$$\begin{aligned} Prob(Y_i = 0|X_i) &= Prob(Y_i^* \leq c|X_i) \\ &= 1 - Prob(\mathcal{E}_i > c - X_i\beta|X_i) \\ &= 1 - Prob(Y_i = 1|X_i) \end{aligned}$$

The expectation (mean) of Y_i conditional on the exogenous variables X_i is:

$$\begin{aligned} E(Y_i|X_i) &= 1 \times Prob(Y_i = 1|X_i) + 0 \times Prob(Y_i = 0|X_i) \\ &= Prob(Y_i = 1|X_i) \\ &= Prob(\mathcal{E}_i > c - X_i\beta|X_i) \end{aligned}$$

The general form of the latent regression model is:

$$E(Y_i|X_i) = Prob(Y_i = 1|X_i) = Prob(\mathcal{E}_i > c - X_i\beta|X_i) = G(X_i\beta)$$

We can predict how the probability of participating in the labour force (Y_i) will change as individuals' characteristics X_i change if we choose either the correct functional form of $G(X_i\beta)$ or the appropriate probability distribution of \mathcal{E}_i .

Note

The **cumulative distribution function(CDF)** of a random variable X is defined as:

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x f_X(t)dt$$

where $f_X(x)$ is the probability density function(PDF) of X. If the functional form of f_X is known, then we can evaluate $P(X \leq x)$ by integration. The PDF can be obtained from the CDF as:

$$\frac{d(f_X(x))}{dx} = f_X(x)$$

Now, if we have a latent regression model of the form:

$$Prob(Y_i = 1|X_i) = Prob(\mathcal{E}_i < c - X_i\beta|X_i)$$

and if the functional form or the probability distribution of \mathcal{E}_i , say, $f_{\mathcal{E}}$ is known, then the response probability can be evaluated since:

$$Prob(\mathcal{E}_i < c - X_i\beta|X_i) = F_{\mathcal{E}}(X_i\beta) = \int_{-\infty}^{X_i\beta} f_{\mathcal{E}}(t)dt$$

The probit and logit models

Setting $G(X_i\beta)$ to be the normal distribution or assuming that \mathcal{E}_i follows the **normal distribution** (that is, $f_{\mathcal{E}}$ is the normal distribution) gives rise to the **probit model**. Since the normal distribution is symmetric about its mean, we have:

$$Prob(\mathcal{E}_i > -X_i\beta|X_i) = Prob(\mathcal{E}_i < X_i\beta|X_i)$$

Thus, the probability of participating in the labour force ($Y_i = 1$) given an individual's

characteristics X_i is given by:

$$\begin{aligned} Prob(Y_i = 1|X_i) &= Prob(\mathcal{E}_i < X_i\beta|X_i) \\ &= \int_{-\infty}^{X_i\beta} f_{\mathcal{E}}(t)dt = \int_{-\infty}^{X_i\beta} \phi(t)dt = \Phi(X_i\beta) \end{aligned}$$

where $\phi(\cdot)$ is the standard normal PDF and $\Phi(\cdot)$ is the standard normal CDF. Thus, we can read the probabilities from the standard normal distribution table.

Setting $G(X_i\beta)$ to be the **logistic distribution** or assuming that \mathcal{E}_i follows the

logistic distribution gives rise to the **logit model**. The logistic distribution function is given by:

$$Prob(\mathcal{E}_i < X_i\beta) = \Lambda(X_i\beta) = \frac{e^{X_i\beta}}{1 + e^{X_i\beta}}$$

Here the response probability $Prob(Y_i = 1|X_i)$ is evaluated as:

$$\begin{aligned} P_i &= Prob(Y_i = 1|X_i) = Prob(\mathcal{E}_i > -X_i\beta|X_i) \quad \text{symmetry} \\ &= 1 - Prob(\mathcal{E}_i < -X_i\beta|X_i) \\ &= 1 - \Lambda(-X_i\beta) \\ &= 1 - \frac{e^{-X_i\beta}}{1 + e^{-X_i\beta}} \\ &= \frac{e^{X_i\beta}}{1 + e^{X_i\beta}} \end{aligned}$$

Similarly, the non-response probability is evaluated as:

$$\begin{aligned} 1 - P_i &= Prob(Y_i = 0|X_i) \\ &= 1 - \frac{e^{X_i\beta}}{1 + e^{X_i\beta}} = \frac{1}{1 + e^{X_i\beta}} \end{aligned}$$

Note that the response and non-response probabilities both lie in the interval $[0, 1]$, and hence, are interpretable.

For the logit model, the ratio:

$$\begin{aligned}\frac{P_i}{1 - P_i} &= \frac{Prob(Y_i = 1|X_i)}{Prob(Y_i = 0|X_i)} = \frac{\frac{e^{X_i\beta}}{1+e^{X_i\beta}}}{\frac{1}{1+e^{X_i\beta}}} \\ &= e^{X_i\beta} = e^{X_{1i}\beta + X_{2i}\beta_2 + \dots + X_{ki}\beta_k}\end{aligned}$$

is the **ratio of the odds** of $Y_i = 1$ against $Y_i = 0$. The natural logarithm of the odds (**log-odds**) is:

$$\ln \left[\frac{P_i}{1 - P_i} \right] = X_i\beta = X_{1i}\beta + X_{2i}\beta_2 + \dots + X_{ki}\beta_k$$

Thus, the log-odds is a linear function of the explanatory variables. The above transformation has certainly helped the popularity of the logit model. Note that for the linear probability model it is P_i that is assumed to be a linear function of the explanatory variables.

Example : Suppose $P_i = Prob(Y_i = 1/X_i)$ is the probability that the i^{th} individual chooses to work given his/her individual characteristics X_i and suppose that the odds is calculated to be $e^{X_i\beta}$. The interpretation of this is that the probability of joining the labour force is **twice as likely** as staying at home given the individual characteristics X_i .

3.4 Estimation of Intrinsically non-linear Models

3.4.1 The maximum likelihood estimation

Bernoulli distribution

Let Y_i be a random variable which can take on the value 1 with probability P_i (called probability of success), and the value 0 with probability $(1 - P_i)$ (called probability of failure). If the observations are independent, then the probability distribution of Y_i is given by:

$$Prob(Y_i = y_i) = P_i^{y_i}(1 - P_i)^{1-y_i} \quad i = 1, 2, \dots, n \quad (3.85)$$

We have seen earlier that:

$$Prob(Y_i = 1|X_i) = G(X_i\beta) \text{ and } Prob(Y_i = 0|X_i) = 1 - G(X_i\beta)$$

Thus, each observation Y_i may be treated as a single draw from a Bernoulli distribution with probability of success ($Y_i = 1$) equal to $G(X_i\beta)$ and probability of failure ($Y_i = 0$) equal to $1 - G(X_i\beta)$. Using equation (3.85), the probability distribution of Y_i is given by:

$$Prob(Y_i = 1|X_i) = [G(X_i\beta)]^{y_i} [1 - G(X_i\beta)]^{1-y_i} \quad i = 1, 2, \dots, n.$$

Since the observations are independent, the **likelihood function** is simply the product of the individual probabilities of the Y_i 's, $i=1, 2, \dots, n$, that is,

$$\begin{aligned} L &= Prob(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n / X) \\ &= Prob(Y_1 = y_1) \times Prob(Y_2 = y_2) \times \dots \times Prob(Y_n = y_n) \\ &= \prod_{i=1}^n [G(X_i\beta)]^{y_i} [1 - G(X_i\beta)]^{1-y_i} \end{aligned}$$

Taking natural logarithms, we get the **log-likelihood function**:

$$\ln(L) = \sum_{i=1}^n [y_i \ln[G(X_i\beta)] + (1 - y_i) \ln[1 - G(X_i\beta)]]$$

Taking the derivative with respect to β and setting the result to zero we have:

$$\frac{\partial \ln(L)}{\partial \beta} = \sum_{i=1}^n \left[\frac{y_i g_i}{G_i} + (1 - y_i) \frac{-g_i}{1 - G_i} \right] X_i = 0 \quad (3.86)$$

where $G_i = G(X_i\beta)$ and $g_i = \frac{\partial G(X_i\beta)}{\partial \beta}$.

These equations are highly non-linear, and we need to apply numerical methods (numerical optimization methods) to obtain the solutions.

If we consider a model with only one explanatory variable X , and if $G(X_i\beta)$ is the logit model:

$$G(X_i\beta) = \frac{e^{X_i\beta}}{1 + e^{X_i\beta}} = \frac{e^{\alpha + X_i\beta}}{1 + e^{\alpha + X_i\beta}}$$

then log likelihood function is given by:

$$\begin{aligned}
\ln(L) &= \sum_{i=1}^n \left[y_i \ln \left(\frac{e^{\alpha+X_i\beta}}{1+e^{\alpha+X_i\beta}} \right) + (1-y_i) \ln \left(1 - \frac{e^{\alpha+X_i\beta}}{1+e^{\alpha+X_i\beta}} \right) \right] \\
&= \sum_{i=1}^n \left[y_i \ln \left(\frac{e^{\alpha+X_i\beta}}{1+e^{\alpha+X_i\beta}} \right) + (1-y_i) \ln \left(\frac{1}{1+e^{\alpha+X_i\beta}} \right) \right] \\
&= \sum_{i=1}^n \left[y_i \ln[e^{\alpha+X_i\beta}] - y_i \ln(1+e^{\alpha+X_i\beta}) - (1-y_i) \ln(1+e^{\alpha+X_i\beta}) \right] \\
&= \sum_{i=1}^n \left[y_i \ln[e^{\alpha+X_i\beta}] - \ln(1+e^{\alpha+X_i\beta}) \right] \\
&= \sum_{i=1}^n \left[y_i \ln(\alpha + X_i\beta) - \ln(1+e^{\alpha+X_i\beta}) \right]
\end{aligned}$$

The first order conditions are:

$$\frac{\partial \ln(L)}{\partial \beta} = 0 \Rightarrow \sum_{i=1}^n \left[y_i - \frac{e^{\tilde{\alpha}+\tilde{\beta}X_i}}{1+e^{\tilde{\alpha}+\tilde{\beta}X_i}} \right] X_i = 0 \quad (3.87)$$

$$\frac{\partial \ln(L)}{\partial \alpha} = 0 \Rightarrow \sum_{i=1}^n \left[y_i - \frac{e^{\tilde{\alpha}+\tilde{\beta}X_i}}{1+e^{\tilde{\alpha}+\tilde{\beta}X_i}} \right] = 0 \quad (3.88)$$

These two equations can be solved for $\tilde{\beta}$ and $\tilde{\alpha}$. Since both equations are non-linear functions of $\tilde{\beta}$ and $\tilde{\alpha}$, the solutions are obtained using numerical methods.

In the classical linear regression model we have $G(X_i\beta) = \alpha + \beta X_i$ and equations (3.87) and (3.88) become:

$$\begin{aligned}
\sum_{i=1}^n [Y_i - (\hat{\alpha} + \hat{\beta}X_i)]X_i &= 0 \Rightarrow \sum_{i=1}^n Y_i X_i = \hat{\alpha} \sum_{i=1}^n X_i + \hat{\beta} \sum_{i=1}^n X_i^2 \\
\sum_{i=1}^n [Y_i - (\hat{\alpha} + \hat{\beta}X_i)] &= 0 \Rightarrow \sum_{i=1}^n Y_i = n\hat{\alpha} + \hat{\beta} \sum_{i=1}^n X_i
\end{aligned}$$

These two equations are simply the **normal equations** which can easily be solved to get the ordinary least squares (OLS) estimators of β and α .

In some non-linear problem, it is most convenient to write down the normal equations and develop a direct (method 1 below) iterative technique for solving them. Whether this works satisfactorily or not – depends on the form of the equations and the iterative method used. There are some of the alternative approaches

- I) Direct search (Trial-and-Error or Derivative-free technique)
- II) Linearization (iterative method or Gauss-Newton Method)
- III) Steepest descent (Direct Optimization)

Example

4 Simultaneous Equations Models (SEM)

4.1 Introduction to SEM

In the previous lessons, we were concerned exclusively with single-equation models, i.e., models in which there was a single dependent variable Y and one or more explanatory variables, the X 's. In such models the emphasis was on estimating and/or predicting the average value of Y conditional upon the fixed values of the X variables. The cause-and-effect relationship, if any, in such models therefore ran from the X 's to the Y (i.e., unidirectional). However, there are situations where there is a two-way flow of influence among economic variables. This occurs if Y is determined by the X 's, and some of the X 's are, in turn, determined by Y . In short, there is a two-way, or simultaneous, relationship between Y and (some of) the X 's, which makes the distinction between dependent and explanatory variables of dubious value.

Simultaneous equations models also differ from most of the econometric models we have considered so far, because they consist of a set of equations. For example, price and quantity are determined by the interaction of two equations, one for supply and the other for demand.

Simultaneous equations models, which contain more than one dependent variable and more than one equation, require special statistical treatment. The least squares estimation procedure is not appropriate in these models, and we must develop new ways to obtain reliable estimates of economic parameters.

Example : At the macro level, aggregate consumption expenditure depends on aggregate disposable income; aggregate disposable income depends upon the national income and taxes imposed by the government; national income depends on aggregate consumption expenditure of the economy.

Disregarding these sequences of relationship, if we estimate a single equation of, say, aggregate consumption on disposable income, then the estimates will be biased and inconsistent.

Example: A simple model of the market for a given commodity may involve a supply and demand

function:

$$Q_t = \alpha_1 + \alpha_2 P_t + \alpha_3 Y_t + U_{1t} \quad (\text{demand}) \quad (4.89)$$

$$Q_t = \beta_1 + \beta_2 P_t + U_{2t} \quad (\text{supply}) \quad (4.90)$$

where Q is the equilibrium quantity exchanged on the market, P is equilibrium price, Y is income of consumers, and u_{1t} and u_{2t} are the disturbance terms. We also have $EU_{1t}^2 = \alpha_1^2$, $EU_{2t}^2 = \alpha_2^2$ and $EU_{1t}U_{2t} = \alpha_{12}$.

Suppose we are interested in the effect of P on Q . Can we toss out the second equation and estimate the first equation alone using OLS?

1. The equilibrium price and quantity are determined in the market by the intersection of supply and demand curves. Therefore, we can not determine equilibrium price by solving the demand equation independently.
2. A shift in the demand function produces a change in both equilibrium price and quantity if the supply curve has an upward slope.

Equations ((4.89),(4.90)) are called the **structural form** of the model under study. These equations can be solved for the ‘endogenous’ variables to give:

$$Q_t = \left(\frac{\alpha_2\beta_1 - \alpha_1\beta_2}{\alpha_2 - \beta_2} \right) - \left(\frac{\alpha_3\beta_2}{\alpha_2 - \beta_2} \right) Y_t + \left(\frac{-\beta_2 U_{1t} + \alpha_2 U_{2t}}{\alpha_2 - \beta_2} \right) \quad (4.91)$$

$$P_t = \left(\frac{-\alpha_1 + \beta_1}{\alpha_2 - \beta_2} \right) - \left(\frac{\alpha_3}{\alpha_2 - \beta_2} \right) Y_t + \left(\frac{-U_{1t} + U_{2t}}{\alpha_2 - \beta_2} \right) \quad (4.92)$$

The solution given by equations (4.91) and (4.92) is called the **reduced form** of the model. The reduced form equations show explicitly how the “endogenous” variables are **jointly dependent** on the “predetermined” variables and the disturbances of the system.

Now from equation (4.92) we have:

$$\begin{aligned}
 E(P_t U_{1t}) &= \left(\frac{-\alpha_1 + \beta_1}{\alpha_2 - \beta_2} \right) E(U_{1t}) - \left(\frac{\alpha_3}{\alpha_2 - \beta_2} \right) Y_t E(U_{1t}) + \left(\frac{-E(U_{1t}^2) + E(U_{1t} U_{2t})}{\alpha_2 - \beta_2} \right) \\
 &= \left(\frac{-\alpha_1 + \beta_1}{\alpha_2 - \beta_2} \right) (0) - \left(\frac{\alpha_3}{\alpha_2 - \beta_2} \right) Y_t (0) + \left(\frac{-\sigma_1^2 + \sigma_{12}}{\alpha_2 - \beta_2} \right) \\
 &= \left(\frac{-\sigma_1^2 + \sigma_{12}}{\alpha_2 - \beta_2} \right) \neq 0
 \end{aligned}$$

Similarly, it can be shown that:

$$E(P_t U_{2t}) = \left(\frac{-\sigma_{12} + \sigma_2^2}{\alpha_2 - \beta_2} \right) \neq 0$$

Thus, in the demand equation (4.89): $Q_t = \alpha_1 + \alpha_2 P_t + \alpha_3 Y_t + U_{1t}$, the variable P_t that appears as an independent or ‘exogenous’ variable is correlated with the disturbance term u_{1t} , and consequently, estimation of the demand equation using OLS leads to **biased** and **inconsistent estimators** of the parameters (refer to **section 2.6** for more details). This is referred to as simultaneity bias.

The solution is to bring the supply function into the picture and estimate the supply and demand functions simultaneously. Such models are known as **simultaneous equations models**.

Example : Wage-price model

$$W_t = \alpha_0 + \alpha_1 U_t + \alpha_2 P_t + U_{1t} \quad (\text{wage equation}) \quad (4.93)$$

$$P_t = \beta_0 + \beta_1 W_t + \beta_2 R_t + \beta_3 M_t + U_{2t} \quad (\text{price equation}) \quad (4.94)$$

where W is rate of change in money wage, U is unemployment rate (in percentage), P is rate of change in prices, R is rate of change in cost of capital, and M is money supply. Here the price variable P enters into the wage equation (4.93) and the wage variable W enters into the price equation (4.94). Thus, these two variables are jointly dependent to each other, and estimation of the two equations individually by OLS yields biased and inconsistent estimators.

Note :

1. **Endogenous variables** are variables that are jointly determined by the economic model. (or are

determined by the exogenous variables).

2. **Exogenous variables** are determined outside of the model and independently of the endogenous variables.

3. **Predetermined variables** are exogenous variables, lagged exogenous variables and lagged endogenous variables. Predetermined variables are non-stochastic and hence independent of the disturbance terms.

4.1.1 Structural form and reduced form of simultaneous equations model (SEM)

Consider the simple Keynesian model of income determination:

$$C_t = \beta_0 + \beta_1 P_t + u_t \quad 0 < \beta_1 < 1 \quad (\text{Consumption function}) \quad (4.95)$$

$$Y_t = C_t + I_t \quad (\text{Income identity}) \quad (4.96)$$

where C is consumption expenditure, Y is income, I is investment (assumed to be exogenous). The above model is said to be the structural form of the SEM, and the parameters β_0 and β_1 are said to be structural parameters. Substituting (4.95) in place of C in (4.96) we get:

$$\begin{aligned} Y_t &= \beta_0 + \beta_1 Y_t + I_t + u_t \\ \Rightarrow Y_t &= \frac{\beta_0}{1 - \beta_1} + \frac{1}{1 - \beta_1} I_t + \frac{1}{1 - \beta_1} u_t \\ &\Rightarrow Y_t = \pi_0 + \pi_1 I_t + w_t \end{aligned} \quad (4.97)$$

Note that equation (4.97) is expressed solely as a function of the exogenous variable I_t and the disturbance term. It is referred to as the reduced form of the SEM, and the parameters π_0 and π_1 are said to be reduced form parameters. Note that the exogenous variable I_t is not correlated with the disturbance term, and hence, we can apply OLS to the reduced form equation to obtain consistent estimators of π_0 and π_1 .

In general, the structural form of a simultaneous system of equations can be described as:

$$\begin{aligned}
\beta_{11}Y_{1t} + \beta_{12}Y_{2t} + \dots + \beta_{1G}Y_{Gt} + \gamma_{11}X_{1t} + \gamma_{12}X_{2t} + \dots + \gamma_{1k}X_{kt} &= u_{1t} \\
\beta_{21}Y_{1t} + \beta_{22}Y_{2t} + \dots + \beta_{2G}Y_{Gt} + \gamma_{21}X_{1t} + \gamma_{22}X_{2t} + \dots + \gamma_{2k}X_{kt} &= u_{2t} \\
\vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots & \\
\beta_{G1}Y_{1t} + \beta_{G2}Y_{2t} + \dots + \beta_{GG}Y_{Gt} + \gamma_{G1}X_{1t} + \gamma_{G2}X_{2t} + \dots + \gamma_{Gk}X_{kt} &= u_{Gt}
\end{aligned}$$

where the Y's are endogenous variables, X's are predetermined variables, and the u's are stochastic disturbances.

- The β 's and γ 's are the structural coefficients
- There are G endogenous and K predetermined variables in the system.
- Not all endogenous and predetermined variables will appear in every equation (that is, some of β 's and γ 's will be zero).
- In each equation, one of the β 's is taken to be unity, that is, one of the endogenous variables serves as the 'dependent' variable when the equation is written out as a standard regression equation.
 - Some of the equations may be identities, that is, their coefficients are known and they contain no stochastic disturbance.

The above model in matrix form is:

$$\begin{bmatrix} \beta_{11} & \beta_{12} & \cdot & \cdot & \cdot & \beta_{1G} \\ \beta_{21} & \beta_{22} & \cdot & \cdot & \cdot & \beta_{2G} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \beta_{G1} & \beta_{G2} & \cdot & \cdot & \cdot & \beta_{GG} \end{bmatrix} \begin{bmatrix} Y_{1t} \\ Y_{2t} \\ \cdot \\ \cdot \\ \cdot \\ Y_{Gt} \end{bmatrix} + \begin{bmatrix} \gamma_{11} & \gamma_{12} & \cdot & \cdot & \cdot & \gamma_{1k} \\ \gamma_{21} & \gamma_{22} & \cdot & \cdot & \cdot & \gamma_{2k} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \gamma_{G1} & \gamma_{G2} & \cdot & \cdot & \cdot & \gamma_{Gk} \end{bmatrix} \begin{bmatrix} X_{1t} \\ X_{2t} \\ \cdot \\ \cdot \\ \cdot \\ X_{kt} \end{bmatrix} = \begin{bmatrix} u_{1t} \\ u_{2t} \\ \cdot \\ \cdot \\ \cdot \\ u_{Gt} \end{bmatrix}$$

$$\Rightarrow BY_t + \Gamma X_t = U_t$$

$$\Rightarrow -B^{-1}\Gamma X_t + B^{-1}U_t \quad (4.98)$$

The reduced form of the system is obtained by expressing the Y's solely as a function of the predetermined variables X's and the disturbance terms:

$$\begin{aligned} Y_{1t} &= \pi_{11}X_{1t} + \pi_{12}X_{2t} + \dots\pi_{1k}X_{kt} + v_{1t} \\ Y_{2t} &= \pi_{21}X_{1t} + \pi_{22}X_{2t} + \dots\pi_{2k}X_{kt} + v_{2t} \\ \vdots & \quad \quad \quad \vdots \quad \quad \quad \vdots \\ Y_{Gt} &= \pi_{G1}X_{1t} + \pi_{G2}X_{2t} + \dots\pi_{Gk}X_{kt} + v_{Gt} \end{aligned}$$

The above model in matrix form is:

$$\underbrace{\begin{bmatrix} Y_{1t} \\ Y_{2t} \\ \vdots \\ Y_{Gt} \end{bmatrix}}_{=Y_t} = \underbrace{\begin{bmatrix} \pi_{11} & \pi_{12} & \cdots & \pi_{1k} \\ \pi_{21} & \pi_{22} & \cdots & \pi_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ \pi_{G1} & \pi_{G2} & \cdots & \pi_{Gk} \end{bmatrix}}_{=\Pi} \underbrace{\begin{bmatrix} X_{1t} \\ X_{2t} \\ \vdots \\ X_{Gt} \end{bmatrix}}_{X_t} + \underbrace{\begin{bmatrix} v_{1t} \\ v_{2t} \\ \vdots \\ v_{Gt} \end{bmatrix}}_{V_t}$$

$$\Rightarrow Y_t = \Pi X_t + V_t \quad (4.99)$$

Comparing (4.98) and (4.99), the relationship between the structural and reduced form parameters is:

$$\Pi = -B^{-1}\Gamma \text{ and } V_t = B^{-1}U_t$$

Example: Suppose we have the following system of equations:

$$q_t = a_1 + b_2p_t + c_1y_t + d_1R_t + U_{1t} \quad (\text{demand function})$$

$$q_t = a_2 + b_2p_t + U_{2t} \quad (\text{supply function})$$

where q is equilibrium quantity exchanged on the market, p is equilibrium price, y is income of

consumers, R is the amount of rainfall (Note: rainfall affects demand, i.e., if there is rain, people do not go shopping) and u_1 and u_2 are the error terms. This structural form model can be re-written as:

$$q_t - b_1 - p_t - a_1 - c_1 y_t - d_1 R_t = u_{1t}$$

$$q_t - b_2 - p_t - a_2 - (0)y_t - (0)R_t = u_{2t}$$

$$\Rightarrow \begin{bmatrix} 1 & -b_1 \\ 1 & -b_2 \end{bmatrix} \begin{bmatrix} q_t \\ p_t \end{bmatrix} + \begin{bmatrix} -a_1 & -c_1 & -d_1 \\ -a_2 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ y_t \\ R_t \end{bmatrix} = \begin{bmatrix} u_{1t} \\ u_{2t} \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} q_t \\ p_t \end{bmatrix} = - \begin{bmatrix} 1 & -b_1 \\ 1 & -b_2 \end{bmatrix}^{-1} \begin{bmatrix} -a_1 & -c_1 & -d_1 \\ -a_2 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ y_t \\ R_t \end{bmatrix} + \begin{bmatrix} 1 & -b_1 \\ 1 & -b_2 \end{bmatrix}^{-1} \begin{bmatrix} u_{1t} \\ u_{2t} \end{bmatrix} \dots\dots\dots (4.1.a)$$

The reduced form equations can be written as:

$$q_t = \pi_1 + \pi_2 y_t + \pi_3 R_t + v_{1t}$$

$$p_t = \pi_4 + \pi_5 y_t + \pi_6 R_t + v_{2t}$$

$$\Rightarrow \begin{bmatrix} q_t \\ p_t \end{bmatrix} = \begin{bmatrix} \pi_1 & \pi_2 & \pi_3 \\ \pi_4 & \pi_5 & \pi_6 \end{bmatrix} \begin{bmatrix} 1 \\ y_t \\ R_t \end{bmatrix} + \begin{bmatrix} v_{1t} \\ v_{2t} \end{bmatrix} \dots\dots\dots (4.1.b)$$

Comparing (4.1.a) and (4.1.b), we can see that:

$$\begin{bmatrix} \pi_1 & \pi_2 & \pi_3 \\ \pi_4 & \pi_5 & \pi_6 \end{bmatrix} = - \begin{bmatrix} 1 & -b_1 \\ 1 & -b_2 \end{bmatrix}^{-1} \begin{bmatrix} -a_1 & -c_1 & -d_1 \\ -a_2 & 0 & 0 \end{bmatrix}$$

$$=$$

$$\frac{1}{b_2 - b_1} \begin{bmatrix} -b_2 & b_1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} -a_1 & -c_1 & -d_1 \\ -a_2 & 0 & 0 \end{bmatrix}$$

$$= \frac{1}{b_2 - b_1} \begin{bmatrix} a_1 b_2 - b_1 a_2 & c_1 b_2 & d_1 b_2 \\ a_1 - a_2 & c_1 & d_1 \end{bmatrix}$$

Thus,

$$\begin{aligned}\pi_1 &= \frac{a_1 b_2 - b_1 a_2}{b_2 - b_1}, \pi_2 = \frac{c_1 b_2}{b_2 - b_1}, \\ \pi_3 &= \frac{d_1 b_2}{b_2 - b_1}, \pi_4 = \frac{a_1 - a_2}{b_2 - b_1}, \\ \pi_5 &= \frac{c_1}{b_2 - b_1}, \text{ and} \\ \pi_6 &= \frac{d_1}{b_2 - b_1}\end{aligned}$$

Note that since the reduced form of the system is obtained by expressing the endogenous variables solely as a function of the predetermined variables, OLS yields consistent estimators of the reduced form parameters. Thus, the OLS estimators $\hat{\pi}_1, \hat{\pi}_2, \hat{\pi}_3, \hat{\pi}_4, \hat{\pi}_5, \hat{\pi}_6$ from equation (4.1.b) above are unbiased and consistent.

4.2 The identification problem

Consider the supply and demand equations in the example above. We stated that the parameters of the reduced form model (that is, the π 's) can be estimated using OLS consistently.

Question: Can we always recover the parameters of the structural equations (that is, a_1, b_1, c_1, a_2, b_2) uniquely from the π 's? In other words, can we always estimate the structural coefficients via the reduced form coefficients? This leads us to the concept of identification. Identification is a problem of model formulation rather than of model estimation or appraisal. We say a model is identified if it is in a unique statistical form, enabling unique estimators of its parameters to be subsequently made from sample data.

Example

In the supply and demand equations in the example above, we can see that:

$$\begin{aligned}\pi_2 &= \frac{c_1 b_2}{b_2 - b_1} \text{ and } \pi_5 = \frac{c_1}{b_2 - b_1} \Rightarrow b_2 = \frac{\pi_2}{\pi_5} \\ \pi_3 &= \frac{d_1 b_2}{b_2 - b_1} \text{ and } \pi_6 = \frac{d_1}{b_2 - b_1} \Rightarrow b_2 = \frac{\pi_3}{\pi_6}\end{aligned}$$

That is, the slope of the supply equation $q_t = a_2 + b_2 p_t + u_{2t}$ **can not be uniquely estimate**. Thus, the model is not identified. This is what we call an **identification problem**.

Note: (Status of identification)

In econometric theory, three possible situations of identifiability can arise: equation under consideration is exactly identified, over identified or under identified.

- 1) If there is a one to one correspondence between the reduced form and structural form parameters, then we have exact identification, that is, there is a unique solution for the structural parameters in terms of the reduced form parameters.
- 2) If the number of reduced form parameters exceeds the number of structural parameters, then we have over identification (no unique solution). Here there is more than sufficient information regarding the equation under consideration.
- 3) If the number of reduced form parameters is less than the number of structural parameters, then we have under identification (no solution). Here there is no sufficient information regarding the equation under consideration.

4.3 Conditions for identification

4.3.1 The order condition for identification

Let G be the total number of endogenous variables in the system and let k be the total number of variables (both endogenous and predetermined) missing from the equation under consideration.

Then if:

- a) $k = G - 1$, the equation is exactly identified.
- b) $k > G - 1$, the equation is over identified.
- c) $k < G - 1$, the equation is under identified

This is known as the order condition for identification. It is a necessary but not sufficient condition for the identification status of an equation.

Example

Wage-price model

$$- W_t = \alpha_0 + \alpha_1 U_t + \alpha_2 P_t + u_{1t} \quad (\text{wage equation})$$

$$- P_t = \beta_0 + \beta_1 W_t + \beta_2 R_t + \beta_3 M_t + u_{2t} \quad (\text{price equation})$$

Here U , R and M are predetermined while W and P are endogenous variables. Thus $G = 2$.

a) Consider the wage equation. The variables R and M are missing from this equation. Thus, $k = 2$. The equation is over identified since $k = 2 > 1 = G^{\sim}1$.

b) Consider the price equation. The variable U is missing from this equation. Thus, $k=1$. The equation is exactly identified since $k = 1 = G - 1$.

Example

The following represents a highly simplified model of the economy:

$$C_t = \alpha_0 + \alpha_1 Y_t + \alpha_2 C_{t-1} + u_{1t} \quad (\text{consumption})$$

$$I_t = \beta_0 + \beta_1 r_t + \beta_2 I_{t-1} + u_{2t} \quad (\text{investment})$$

$$r_t = \gamma_0 + \gamma_1 Y_t + \gamma_2 M_t + u_{3t} \quad (\text{money market})$$

$$Y_t = C_t + I_t + G_t \quad (\text{income identity})$$

where C = consumption, Y = income, I = investment, r = rate of interest, M = money supply and G = government expenditure. The variables C_t, I_t, Y_t, r_t are endogenous while the remaining are predetermined variables. Thus, $G = 4$.

a) Consider the consumption equation. The variables $I_{t-1}, I_t, M_t, r_t, G_t$ are missing from this equation. Thus, $k = 5$. The equation is over identified since $k = 5 > G^{\sim}1 = 3$.

b) Consider the investment equation. The variables $C_{t-1}, C_t, M_t, Y_t, G_t$ are missing from this equation. Thus, $k = 5$. The equation is over identified since $k = 5 > G^{\sim}1 = 3$.

4.3.2 The rank condition for identification

The rank condition states that: in a system of G equations any particular equation is (exactly or over) identified if and only if it is possible to construct at least one non-zero determinant of order (G-1) from the coefficients of the variables excluded from that particular equation but contained in the other equations of the system.

Example

Consider the simplified model of the economy above:

$$C_t = \alpha_0 + \alpha_1 Y_t + \alpha_2 C_{t-1} + u_{1t} \quad (\text{consumption})$$

$$I_t = \beta_0 + \beta_1 r_t + \beta_2 I_{t-1} + u_{2t} \quad (\text{investment})$$

$$r_t = \gamma_0 + \gamma_1 Y_t + \gamma_2 M_t + u_{3t} \quad (\text{money market})$$

$$Y_t = C_t + I_t + G_t \quad (\text{income identity})$$

This model may be re-written as:

$$-C_t + \alpha_0 + \alpha_1 Y_t + \alpha_2 C_{t-1} + 0I_t + 0r_t + 0I_{t-1} + 0M_t + 0G_t + u_{1t} = 0$$

$$0C_t + \beta_0 + 0Y_t + 0C_{t-1} - I_t + \beta_1 r_t + \beta_2 I_{t-1} + 0M_t + 0G_t + u_{2t} = 0$$

$$0C_t + \gamma_0 + \gamma_1 Y_t + 0C_{t-1} + 0I_t - r_t + 0I_{t-1} + \gamma_2 M_t + 0G_t + u_{3t} = 0$$

$$C_t + 0 - Y_t + 0C_{t-1} + I_t + 0r_t + 0I_{t-1} + 0M_t + G_t + 0 = 0$$

Note that the coefficient of a variable excluded from an equation is equal to zero. Ignoring the random disturbances and the constants, a table of the parameters of the model is as follows:

	C_t	Y_t	C_{t-1}	I_t	r_t	I_{t-1}	M_t	G_t
consumption	-1	α_1	α_2	0	0	0	0	0
investment	0	0	0	-1	β_1	β_2	0	0
money market	0	γ_1	0	0	-1	0	γ_2	0
income identity	1	-1	0	1	0	0	0	1

Now suppose we want to check the identification status of the consumption function.

- We eliminate the row corresponding to the consumption function.
- We eliminate the columns in which the consumption function has non-zero coefficients.

The two steps are shown below:

	C_t	Y_t	C_{t-1}	I_t	r_t	I_{t-1}	M_t	G_t
consumption	-1	α_1	α_2	0	0	0	0	0
investment	0	0	0	-1	β_1	β_2	0	0
money market	0	γ_1	0	0	-1	0	γ_2	0
income identity	1	-1	0	1	0	0	0	1

Note that by doing steps (a) and (b) above, we are left with the coefficients of variables not included in the consumption function, but contained in the other equations of the system.

After eliminating the relevant row and columns, we get the following table (matrix) of parameters:

I_t	r_t	I_{t-1}	M_t	G_t
-1	β_1	β_2	0	0
0	-1	0	γ_2	0
1	0	0	0	1

.....(*)

Since the system has $G = 4$ equations, form the determinants of order $(G-1) = 3$ and examine their value.

- If at least one of these determinants is non-zero, then the consumption equation is (exactly or over) identified.
- If all determinants of order 3 are zero, then the consumption equation is under identified.

For example;

$$\Delta_1 = \begin{vmatrix} -1 & \beta_1 & \beta_2 \\ 0 & -1 & 0 \\ 1 & 0 & 0 \end{vmatrix} = - \begin{vmatrix} -1 & 0 \\ 0 & 0 \end{vmatrix} - \beta_1 \begin{vmatrix} 1 & 0 \\ 0 & 0 \end{vmatrix} + \beta_2 \begin{vmatrix} 0 & -1 \\ 1 & 0 \end{vmatrix} = -1(0) + \beta_1(0) + \beta_2(1) = \beta_2 \neq 0$$

or

$$\Delta_2 = \begin{vmatrix} \beta_2 & 0 & 0 \\ 0 & \gamma_2 & 0 \\ 0 & 0 & 1 \end{vmatrix} = \beta_2(\gamma_2) \neq 0$$

Thus, we can form at least one non-zero determinant of order 3, and hence, the consumption equation is exactly or over identified.

To see whether the consumption equation is exactly or over identified, we can use the order condition. Since we have four endogenous variables (C_t, I_t, Y_t, r_t) , $G = 4$. As can be seen from Table (*) above, the variables $I_{t-1}, I_t, M_t, r_t, G_t$ are missing from the consumption equation, meaning $k = 5$. Thus, the equation is over identified since $k = 5 > G - 1 = 3$.

4.4 Estimation of simultaneous equations models: ILS, 2SLS

4.4.1 Indirect least squares (ILS) method

In this method, we first obtain the estimates of the reduced form parameters by applying OLS to the reduced form equations and then indirectly get the estimates of the parameters of the structural model. This method is applied to exactly identified equations.

Steps:

- a) Obtain the reduced form equations (that is, express the endogenous variables in terms of predetermined variables).
- b) Apply OLS to the reduced form equations individually. OLS will yield consistent estimates of the reduced form parameters (since each equation involves only non-stochastic (predetermined) variables that appear as ‘independent’ variables).
- c) Obtain (or recover back) the estimates of the original structural coefficients using the estimates in step (2).

Example

Consider the following model for demand and supply of pork:

$$Q_t = a_1 + a_2 P_t + a_3 Y_t + u_{1t} \quad (\text{demand function}) \quad (4.100)$$

$$Q_t = a_2 + a_3 P_t + a_3 Z_t + u_{2t} \quad (\text{supply function}) \quad (4.101)$$

where Q_t is consumption of pork (pounds per capita), P_t is real price of pork (cents per pound), Y_t is disposable personal income (dollars per capita) and Z_t is ‘predetermined elements in pork production’.

Here P and Q are endogenous variables while Y and Z are predetermined variables. It can easily be shown that both equations are exactly identified. Thus, we can apply ILS to estimate the parameters. We first express P and Q in terms of the predetermined variables and disturbances as:

$$Q_t = \left(\frac{b_2 a_1 - b_1 a_2}{b_2 - b_1} \right) + \left(\frac{c_1 b_2}{b_2 - b_1} \right) Y_t - \left(\frac{c_2 b_1}{b_2 - b_1} \right) Z_t + \left(\frac{b_2 u_{1t} - b_1 u_{2t}}{b_2 - b_1} \right) \quad (4.102)$$

$$P_t = \left(\frac{a_1 - a_2}{b_2 - b_1} \right) + \left(\frac{c_1}{b_2 - b_1} \right) Y_t - \left(\frac{c_2}{b_2 - b_1} \right) Z_t + \left(\frac{u_{1t} - u_{2t}}{b_2 - b_1} \right) \quad (4.103)$$

We can re-write equations (4.102) and (4.103) as:

$$q_t = \pi_1 + \pi_2 Y_t + \pi_3 Z_t + \epsilon_{1t}$$

$$p_t = \pi_4 + \pi_5 Y_t + \pi_6 Z_t + \epsilon_{1t}$$

$$\frac{\pi_2}{\pi_5} = \frac{c_1 b_2 / b_2 - b_1}{c_1 b_2 - b_1} = b_2 \Rightarrow \hat{b}_2 = \frac{\hat{\pi}_2}{\hat{\pi}_5}$$

$$\frac{\pi_3}{\pi_6} = \frac{c_2 b_1 / b_2 - b_1}{c_2 / b_2 - b_1} = b_1 \Rightarrow \hat{b}_1 = \frac{\hat{\pi}_3}{\hat{\pi}_6}$$

$$\pi_5 = \frac{c_1}{b_2 - b_1} \Rightarrow c_1 = \pi_5 (b_2 - b_1) \Rightarrow \hat{c}_1 = \hat{\pi}_5 (\hat{b}_2 - \hat{b}_1)$$

Similarly, it can be shown that $\hat{c}_2 = \hat{\pi}_6 (\hat{b}_2 - \hat{b}_1)$, $\hat{a}_1 = \hat{\pi}_1 - \hat{b}_1 \hat{\pi}_4$ and $\hat{a}_2 = \hat{\pi}_1 - \hat{b}_2 \hat{\pi}_4$

4.4.2 Instrumental variable(IV) method

Suppose we have the model (in deviation form):

$$y_i = \beta x_i + \mathcal{E}_i$$

where x_i is correlated with \mathcal{E}_i . We can not estimate β by OLS as it will yield an inconsistent estimator of β (refer to errors var variables for details). What we do is search for an instrumental variable (IV) z_i that is uncorrelated with \mathcal{E}_i but correlated with x_i ; that is, $cov(z_i, \mathcal{E}_i) = 0$ and $cov(z_i, x_i) \neq 0$. The sample counterpart of $cov(z_i, \mathcal{E}_i) = 0$ is:

$$\begin{aligned} \frac{1}{n} \sum z_i \mathcal{E}_i = 0 &\implies \frac{1}{n} \sum z_i (y_i - \beta x_i) = 0 \\ &\implies \frac{1}{n} \sum z_i y_i = \hat{\beta} \left(\frac{1}{n} \sum z_i x_i \right) \\ &\implies \hat{\beta} = \frac{\frac{1}{n} \sum z_i y_i}{\frac{1}{n} \sum z_i x_i} = \frac{\sum z_i y_i}{\sum z_i x_i} \end{aligned}$$

$\hat{\beta}$ can be expressed as:

$$\hat{\beta} = \frac{\sum z_i y_i}{\sum z_i x_i} = \frac{\sum z_i (\beta x_i + \mathcal{E}_i)}{\sum z_i x_i} = \beta + \frac{\sum z_i \mathcal{E}_i}{\sum z_i x_i}$$

Now we have,

$$plim(\sum z_i \mathcal{E}_i / n) = cov(z_i, \mathcal{E}_i) = 0$$

$$plim(\sum z_i x_i / n) = cov(z_i, x_i) \neq 0$$

Thus,

$$plim(\hat{\beta}) = \beta + plim\left(\frac{\sum z_i \mathcal{E}_i}{\sum z_i x_i}\right) = \beta + \frac{plim(\sum z_i \mathcal{E}_i / n)}{plim(\sum z_i x_i / n)} = \beta + \frac{0}{\neq 0} = \beta$$

that is, the IV estimator $\hat{\beta}$ is a consistent estimator of β .

Consider the following simultaneous equations model:

$$y_1 = a_1 + b_1 y_2 + c_1 z_1 + c_2 z_2 + u_1$$

$$y_2 = a_2 + b_2 y_1 + c_3 z_3 + u_2$$

where y_1 and y_2 are endogenous while z_1 , z_2 and z_3 are predetermined.

Consider the estimation of the first equation:

- Since z_1 and z_2 are predetermined, they are not correlated with u_1 , that is, $cov(z_1, u_1) = 0$ and $cov(z_2, u_1) = 0$

- y_2 is not independent of u_2 , that is, $cov(y_2, u_1) \neq 0$

Thus, OLS method of estimation can not be applied. To find consistent estimators, we look for a variable that is correlated with y_2 but not correlated with u_1 . Fortunately we have z_3 that satisfies these two conditions, that is, $cov(y_2, z_3) \neq 0$ and $cov(z_3, u_1) = 0$. Thus, z_3 can serve as an IV for y_2 .

The procedure for estimation of the first equation is as follows:

a) Regress y_2 on z_1 , z_2 and z_3 ; that is, using OLS estimate the model:

$$y_2 = a_{10} + a_{11}z_1 + a_{12}z_2 + a_{13}z_3 + v_1.$$

b) Obtain \hat{y}_2 where $\hat{y}_2 = \hat{a}_{10} + \hat{a}_{11}z_1 + \hat{a}_{12}z_2 + \hat{a}_{13}z_3$.

c) Regress y_2 on \hat{y}_2 , z_1 and z_2 ; that is, estimate the model:

$$y_1 = a_1 + b_1\hat{y}_2 + c_1z_1 + c_2z_2 + u_1$$

Note that since z_1 , z_2 and z_3 are predetermined variables, and hence, not correlated with u_1 , we have:

$$\text{cov}(\hat{y}_2, u_1) = \text{cov}(\hat{a}_{10} + \hat{a}_{11}z_1 + \hat{a}_{12}z_2 + \hat{a}_{13}z_3, u_1) = 0$$

Thus, the OLS estimation using the above procedure yields consistent estimators.

Consider the second equation.

- Since z_3 is predetermined, it is not correlated with u_2 , that is, $\text{Cov}(z_3, u_2) = 0$.
- y_1 is not independent of u_2 , that is, $\text{Cov}(y_1, u_2) \neq 0$

Again OLS can not be applied to estimate the parameters. To find consistent estimators, we look for a variable that is correlated with y_1 but not correlated with u_2 . Here we have two choices, namely, z_1 and z_2 that can serve as instruments.

Note : We have more than enough instrumental variables since the second equation is over identified.

In order to estimate the second equation:

a) Regress y_1 on z_1 and z_3 (if z_1 is considered as an IV for y_1) or regress y_1 on z_2 and z_3 (if z_2 is considered as an IV for y_1) using OLS and obtain \hat{y}_1 .

b) Regress y_2 on \hat{y}_1 and z_3 ; that is, estimate the model:

Note that the solution is not unique, that is, depending on whether z_1 is considered as an IV for y_1 or z_2 is considered as an IV for y_1 , we may get different results.

Example

Using a data on some characteristics of the wine industry in Australia, which is saved as “wine.csv”.

```
wine<-read.csv("wine.csv",header = T,sep = ",")
wine_var<-wine[1,] # To look the variables for short codes
names(wine)=c("year","consumption","storage_costs","price_wine",
              "price_beer","advertising","income")
# Converting to log
year<-wine$year; wine=wine[, -1]; wine=log(wine); wine<-cbind(year,wine)
```

It is assumed that a reasonable demand-supply model for the industry would be (where all variables are in logs):

$$Q_1 = a_0 + a_1PW_t + a_2PB_t + a_3Y_t + a_4A_t + u_t \quad \text{.....(demand)}$$

$$Q_1 = b_0 + b_1PW_t + b_2S_t + v_t \quad \text{.....(supply)}$$

where Q is real per capita consumption of wine, PW is the price of wine relative to CPI, PB is the price of beer relative to CPI, Y is real per capita disposable income, A is real per capita advertising expenditure, and S is index of storage costs. Here Q and PW are the two endogenous variables while the rest are exogenous variables.

Apply instrumental variables method of estimation?

solution

Using Instrumental variables method of Estimation

To estimate the demand function we have only one instrumental variable (IV) S. But for the estimation of the supply we have available three IVs: PB, Y and A.

I) Estimation of the supply function

The supply equation is **over-identified**. Thus, we have three possible IV's for price of wine (PW):

- a) price of beer (PB),
- b) advertising expense (A),
- c) income (Y).

(a) Now let's take price of beer as IV to estimate the supply function:

First we regress price_wine(PW) on price_beer(PB) and storage_costs(S), and obtain the predicted values of PW (IV pb for prcwine). At last, we estimate the supply function by regressing

consumption (Q) on S and (IV pb for prcwine).

The results are as follows:

```
rm(list = ls())
wine<-read.csv("wine.csv",header = T,sep = ",")
wine_var<-wine[1,] # To look the variables for short codes
names(wine)=c("year","consumption","storage_costs","price_wine",
              "price_beer","advertising","income")
year<-wine$year; wine=wine[,-1];wine=log(wine);wine<-cbind(year,wine)

# regress price_wine(PW) on price_beer(PB) and storage_costs(S)
prcwine<-lm(price_wine~storage_costs+price_beer,data=wine)
#obtain the predicted values of price_wine(PW)
pred_prcwine_4_pb<-fitted(prcwine)
# estimate the supply function by regressing consumption(Q) on
#pred_prcwine and storage_costs
fit_supply<-lm(consumption~pred_prcwine_4_pb+storage_costs,data=wine)
summary(fit_supply)
```

```
##
## Call:
## lm(formula = consumption ~ pred_prcwine_4_pb + storage_costs,
##     data = wine)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.7108 -0.1949 -0.0336  0.2741  0.4126
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -10.759     44.160   -0.24    0.81
```

```
## pred_prcwine_4_pb      0.336      16.604      0.02      0.98
## storage_costs          2.131       6.876      0.31      0.76
##
## Residual standard error: 0.314 on 17 degrees of freedom
## Multiple R-squared:  0.788,   Adjusted R-squared:  0.764
## F-statistic: 31.7 on 2 and 17 DF,  p-value: 1.84e-06
```

as shown above the Instrumental variable price of beer for price of wine is not significant.

(b) Now let's take advertising expense as IV to estimate the supply function:

First we regress price of wine (PW) on advertising expense (A) and storage cost (S), and obtain the predicted values of PW (IV ad for prcwine). At last, we estimate the supply function by regression consumption (Q) on S and (IV ad for prcwine).

The results are as follows:

```
rm(list = ls())
wine<-read.csv("wine.csv",header = T,sep = ",")
wine_var<-wine[1,] # To look the variables for short codes
names(wine)=c("year","consumption","storage_costs","price_wine",
              "price_beer","advertising","income")
year<-wine$year; wine=wine[,-1];wine=log(wine);wine<-cbind(year,wine)

prcwine_adv<-lm(price_wine~advertising+storage_costs,data=wine)
# the predicted values of price_wine(PW) is
pred_prcwine_4_adv<-fitted(prcwine_adv)
#Then estimate the supply function by regressing consumption(Q)
# on storage_costs and
fit_supply_adv<-lm(consumption~pred_prcwine_4_adv+storage_costs,
                  data=wine)
summary(fit_supply_adv)

##
```

```
## Call:
## lm(formula = consumption ~ pred_prcwine_4_adv + storage_costs,
##     data = wine)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.5200 -0.1409 -0.0741  0.1859  0.5105
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      -17.65       4.62   -3.82  0.0014 **
## pred_prcwine_4_adv    2.93       1.67    1.75  0.0981 .
## storage_costs     1.06       0.74    1.43  0.1709
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.289 on 17 degrees of freedom
## Multiple R-squared:  0.821, Adjusted R-squared:  0.8
## F-statistic: 38.9 on 2 and 17 DF, p-value: 4.51e-07
```

(c) Now let's take income as IV to estimate the supply function:

First we regress price of wine (PW) on income (Y) and storage cost (S), get the predicted values of PW (IV inc for prcwine), and then estimate the supply function by regressing consumption (Q) on S and (IV inc for prcwine).

```
rm(list = ls())
wine<-read.csv("wine.csv",header = T,sep = ",")
wine_var<-wine[1,] # To look the variables for short codes
names(wine)=c("year","consumption","storage_costs","price_wine",
              "price_beer","advertising","income")
year<-wine$year; wine=wine[, -1]; wine=log(wine)
```

```
wine<-cbind(year,wine)

prcwine_inc<-lm(price_wine~income+storage_costs,data=wine)
# the predicted values of price_wine(PW) is
pred_prcwine_4_inc<-fitted(prcwine_inc)
#Then estimate the supply function by regressing consumption(Q)
# on storage_costs and
fit_supply_inc<-lm(consumption~storage_costs+pred_prcwine_4_inc,
                    data=wine)
summary(fit_supply_inc)

##
## Call:
## lm(formula = consumption ~ storage_costs + pred_prcwine_4_inc,
##     data = wine)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.3421 -0.1083  0.0101  0.0827  0.3727
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -16.979     1.343  -12.64  4.5e-10 ***
## storage_costs      1.163     0.234   4.97  0.00012 ***
## pred_prcwine_4_inc  2.676     0.422   6.34  7.5e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.171 on 17 degrees of freedom
## Multiple R-squared:  0.937, Adjusted R-squared:  0.93
```

```
## F-statistic: 127 on 2 and 17 DF, p-value: 6.18e-11
```

By comparing the estimated models using the three IV's, it seems that income is the best IV as the resulting estimated model has the highest coefficient of determination ($R^2 = 0.9371$). Since all variables are in logs, the coefficients are elasticities. Thus, quantity supplied is responsive to both price and storage costs (both p-values < 0.001). In particular, the price elasticity of supply for wine is about 2.68.

II) Estimation of the demand function

The demand equation is **exactly-identified**. Thus, we just have one available IV: storage costs.

First we regress price of wine (PW) on price of beer (PB), income (Y), advertising expense (A) and storage cost (S), get the predicted values of PW (Predicted price of wine), and then estimate the demand function by regressing consumption (Q) on (Predicted price of wine), PB, Y and A.

The results are as follows:

```
rm(list = ls())
wine<-read.csv("wine.csv",header = T,sep = ",")
wine_var<-wine[1,] # To look the variables for short codes
names(wine)=c("year","consumption","storage_costs","price_wine",
              "price_beer","advertising","income")
year<-wine$year; wine=wine[,-1];wine=log(wine);wine<-cbind(year,wine)

prcwine_stors<-lm(price_wine~price_beer+income+advertising+storage_costs,
                  data=wine)

# the predicted values of price_wine(PW) is
predicted_PW<-fitted(prcwine_stors)

#Then to estimate the demand function by regressing consumption(Q)
fit_dem<-lm(consumption~predicted_PW+price_beer+income+advertising,
            data=wine)

summary(fit_dem)
```

```
##
```

```
## Call:
```



```
## lm(formula = consumption ~ predicted_PW + price_beer + income +
##     advertising, data = wine)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.28307 -0.09206 -0.00374  0.08710  0.29299
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -26.195      6.594   -3.97   0.0012 **
## predicted_PW    0.643      0.838    0.77   0.4546
## price_beer    -0.140      0.878   -0.16   0.8757
## income         4.082      1.594    2.56   0.0217 *
## advertising   -0.985      0.835   -1.18   0.2563
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.155 on 15 degrees of freedom
## Multiple R-squared:  0.955, Adjusted R-squared:  0.943
## F-statistic: 79 on 4 and 15 DF, p-value: 6.8e-10
```

We observe that the coefficient of determination is 95.47% and the F-statistic is significant. However, most of the regression coefficients are insignificant. Furthermore, all the coefficients except that of income (Y) have the wrong signs. This is probably due to multicollinearity (MC) and check it using vif function from 'car' package.

```
#install.packages("car")
library(car)
vif(fit_dem)
```

```
## predicted_PW price_beer income advertising
##      11.718      2.551    70.910      32.924
```

This shows that the variance inflation factor (VIF) for income and advertising expense are large (far greater than 10). Thus, we have to drop one of them. From practical point of view, it seems wise to drop advertising expense and re-estimate the model. The results are:

```
#library(car)
fit_dem_adv_drop<-lm(consumption~predicted_PW+price_beer+income,
                      data=wine)
vif(fit_dem_adv_drop)
```

```
## predicted_PW    price_beer      income
##           5.213         2.520         7.847
```

The problem of MC is now solved as the VIF's are greatly reduced (all less than 10).

```
summary(fit_dem_adv_drop)

##
## Call:
## lm(formula = consumption ~ predicted_PW + price_beer + income,
##     data = wine)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.3396 -0.0995  0.0076  0.0682  0.3233
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -21.109      5.052   -4.18  0.00071 ***
## predicted_PW    1.380      0.566    2.44  0.02673 *
## price_beer    -0.252      0.883   -0.29  0.77869
## income         2.308      0.537    4.30  0.00055 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

##

Residual standard error: 0.156 on 16 degrees of freedom

Multiple R-squared: 0.95, Adjusted R-squared: 0.941

F-statistic: 102 on 3 and 16 DF, p-value: 1.18e-10

However, the coefficients of both price of wine and price of beer have wrong signs. In particular, the coefficient of price of wine not only has the wrong sign but is also significant. This is difficult to interpret. For the other variables, the conclusion we arrive at is that the demand for wine is not responsive to the price of beer, but is responsive to income. The income elasticity of demand for wine is about 2.31.

4.4.3 Two-stage least squares (2-SLS) method

The main difference between the IV and 2-SLS methods is that in the former case the \hat{y}_i are used as instruments, while in the latter case the \hat{y}_i are used as regressors. Both methods yield the same result if the equation under consideration is exactly identified. The 2-SLS procedure is generally applicable for estimation of over-identified equations as it provides unique estimators.

Steps :

- a) Estimate the reduced form equations by OLS and obtain the predicted \hat{y}_i .
- b) Replace the right hand side endogenous variables in the structural equations by the corresponding \hat{y}_i and estimate them by OLS.

Consider the above simultaneous equations model:

$$y_1 = a_1 + b_1 y_2 + c_1 z_1 + c_2 z_2 + u_1 \quad (4.104)$$

$$y_2 = a_2 + b_2 y_1 + c_3 z_3 + u_2 \quad (4.105)$$

where y_1 and y_2 are endogenous while z_1 , z_2 and z_3 are predetermined.

Since $cov(y_2, u_1) \neq 0$ and $Cov(y_1, u_2) \neq 0$, we can not apply OLS. Since equation (4.104) is exactly identified, the 2-SLS procedure is the same as the IV method. The 2-SLS procedure of estimation of equation (4.105) (which is over-identified) is: . We first estimate the reduced form equations by OLS; that is, we regress y_1 on z_1 , z_2 and z_3 using OLS and obtain \hat{y}_1 .

. We then replace y_1 by \hat{y}_1 and estimate equation (4.105) by OLS, that is, we apply OLS to:

$$y_2 = b_2\hat{y}_1 + c_3z_3 + u_2$$

Note

- a) Unlike ILS, 2-SLS provides only one estimate per parameter for over-identified models.
- b) In case of exactly identified equations, both ILS and 2-SLS produce the same parameter estimates.

Example

for the previous data, 'wine.csv', apply two stages least squares (2-SLS) method of estimation to the supply-demand function of wine. where the demand-supply model of the industry is (where all variables are in logs):

$$\begin{aligned} Q_1 &= a_0 + a_1PW_t + a_2PB_t + a_3Y_t + a_4A_t + u_t && \text{.....(demand)} \\ Q_1 &= b_0 + b_1PW_t + b_2S_t + v_t && \text{.....(supply)} \end{aligned}$$

where Q is real per capita consumption of wine, PW is the price of wine relative to CPI, PB is the price of beer relative to CPI, Y is real per capita disposable income, A is real per capita advertising expenditure, and S is index of storage costs. Here Q and PW are the two endogenous variables while the rest are exogenous variables.

Apply 2-SLS method of estimation?

```
rm(list = ls())
wine<-read.csv("wine.csv",header = T,sep = ",")
wine_var<-wine[1,] # To look the variables for short codes
names(wine)=c("year","consumption","storage_costs","price_wine",
              "price_beer","advertising","income")
# Converting to log
year<-wine$year; wine=wine[,-1];wine=log(wine)
wine<-cbind(year,wine)
wine_var # Look the variable
```

solution

Estimation using two stages least squares (2-SLS)

In this method, we first find the reduced form equations by regressing each endogenous variable on all exogenous variables. Then we replace all the endogenous variables in each equation by their predicted values from the reduced forms and estimate each equation by OLS. Note that the IV estimator and the 2-SLS estimator are the same if the equation under consideration is exactly identified. In our case we have seen that the demand equation is exactly identified. Thus, the IV and 2-SLS estimators of the parameters are the same.

(I)Estimate supply function

To estimate the supply function: We first regress the price of wine (PW) on all exogenous variables PB, Y, A and S, and get the predicted values (PW 2sls). We then estimate the supply function by regressing consumption (Q) on S and PW 2sls.

The results are shown below:

```
# regress the price_wine on price_beer, income, advertising,
#      and storage_costs.
pw2sls<-lm(price_wine~price_beer+income+advertising+storage_costs,
           data=wine)
# predicted value of price wine .
pred_prcwine<-fitted(pw2sls)
#Estimate the supply function by regressing consumption (Q)
# on pred_prcwine and
fit_sup<-lm(consumption~pred_prcwine+storage_costs,data=wine)
summary(fit_sup)

##
## Call:
## lm(formula = consumption ~ pred_prcwine + storage_costs, data = wine)
##
## Residuals:
##      Min      1Q    Median      3Q      Max
```

```
## -0.28851 -0.09984 -0.00019 0.08531 0.27882
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -16.820      1.089  -15.44 1.9e-11 ***
## pred_prcwine    2.616      0.334   7.83 4.9e-07 ***
## storage_costs  1.188      0.192   6.19 9.9e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.146 on 17 degrees of freedom
## Multiple R-squared:  0.954, Adjusted R-squared:  0.949
## F-statistic: 176 on 2 and 17 DF, p-value: 4.26e-12
```

We can see from the above table that the coefficients of both price of wine and storage cost are significant. The price elasticity of supply is about 2.62.

(II) Estimate demand function

To estimate the demand equation, we have to regress the consumption on all exogenous variables of the equation, i.e., price of wine (but use pred_prcwine) (PW), price of beer (PB), income (Y) and advertising (A).

```
fit_dem<-lm(consumption~pred_prcwine+price_beer+income+advertising,
            data=wine)
summary(fit_dem)
```

```
##
## Call:
## lm(formula = consumption ~ pred_prcwine + price_beer + income +
##     advertising, data = wine)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -0.28307 -0.09206 -0.00374 0.08710 0.29299
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -26.195      6.594   -3.97  0.0012 **
## pred_prcwine    0.643      0.838    0.77  0.4546
## price_beer    -0.140      0.878   -0.16  0.8757
## income         4.082      1.594    2.56  0.0217 *
## advertising   -0.985      0.835   -1.18  0.2563
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.155 on 15 degrees of freedom
## Multiple R-squared:  0.955, Adjusted R-squared:  0.943
## F-statistic: 79 on 4 and 15 DF, p-value: 6.8e-10
```

check for MC problem

```
vif(fit_dem)
```

```
## pred_prcwine price_beer income advertising
##          11.718         2.551        70.910         32.924
```

As shown above there is multicollinearity problem. So, let's drop advertising as we did in the previous.

```
fit_dem_drop_ad<-lm(consumption~pred_prcwine+price_beer+income,
                    data=wine)
summary(fit_dem_drop_ad)

##
## Call:
## lm(formula = consumption ~ pred_prcwine + price_beer + income,
##     data = wine)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.3396 -0.0995  0.0076  0.0682  0.3233
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -21.109      5.052   -4.18  0.00071 ***
## pred_prcwine    1.380      0.566    2.44  0.02673 *
## price_beer    -0.252      0.883   -0.29  0.77869
## income         2.308      0.537    4.30  0.00055 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.156 on 16 degrees of freedom
## Multiple R-squared:  0.95,    Adjusted R-squared:  0.941
## F-statistic: 102 on 3 and 16 DF,  p-value: 1.18e-10

## check MC problem
vif(fit_dem_drop_ad)

## pred_prcwine  price_beer      income
##          5.213          2.520          7.847
```

Since demand equation is **exactly identified**, the 2SLS estimators of the parameters are the same with IV estimators that we have done in the previous.