

Econometric Modelling
Professor Sujata Kar
Department of Management Studies
Indian Institute of Technology, Roorkee
Lecture 15
Omitted Variables and Parameters Stability - I

Hello and welcome back to the course on Econometric Modelling. This is Module 15 and it deals with Omitted Variables and Parameters Stability.

(Refer Slide Time: 00:36)

Part 1: Introduction to Econometrics Module 1: An Overview Module 2: Formulation of Econometric Modelling Module 3 & 4: Review of Basic Concepts Module 5: Types of Data	Part 5: Univariate Time Series Modeling Module 25, 26, 27: Problem of Serial Correlation Module 28: AR, MA & ARMA Processes Module 29: Modelling Seasonal Variations
Part 2: Overview of Classical Linear Regression Model Module 6 & 7: Simple Regression Module 8: Assumption of Classical Linear Regression Module 9: Properties of OLS Estimators Module 10: Hypothesis Testing	Part 6: Models with Binary Dependent and Independent Variables Module 30 & 31: Spline Function & Categorical Variables Module 32 & 33: Probit, Logit and Multinomial Logit Models
Part 3: Multiple Regression Analysis & Diagnostic Tests Module 11, 12 & 13: Multiple Regression Module 14: Problems of Multicollinearity Module 15 & 16: Omitted Variables & Parameter Stability Module 17 & 18: Problem of Heteroscedasticity	Part 7: Multivariate Models Module 33 & 34: Simultaneous Equations System Module 35 & 36: Introduction to VARs
Part 4: Statistical Inference Module 19: t-test Module 20 & 21: Wald test Module 22 & 23: F-test Module 24: Chow test	Part 8: Modelling Long Run Relationships Module 37, 38 & 39: Stationarity & Unit Root Testing Module 40: Basics of Cointegration

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE

Module 15 and 16 are both devoted to the problem of basically model misspecification. So, in model 15, I present two types of model misspecification, omitted variable problems and measurement errors. While in module 16 we will discuss how to deal with such problems that if variables have been omitted or there is a measurement error, then what are the steps required to be taken in order to correct for those kinds of model misspecifications.

(Refer Slide Time: 01:14)

Omitted Variables & Parameter Stability

MODULE - 15

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE 3

Model Misspecification

- Misspecification of models can have serious implications.
- For example, suppose,
- $\ln(wage) = \beta_0 + \beta_1 edu + \beta_2 exp + \beta_3 exp^2 + u$ $E(u|x) \neq 0$
- But if we exclude exp^2 from this specification, then that violates the Gauss-Markov assumptions and we will not get unbiased estimators.
- Similarly, suppose $\ln(wage) = \beta_0 + \beta_1 edu + \beta_2 exp + \beta_3 abil + u$
- However, since ability is not observed and is not included, then we have the problem of misspecification.

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE 4

Misspecification of models can have serious implications. For example, suppose this is our model of the world (*refer slide time:01:14*) that is the logarithm of wages regressed on education experience and experience square plus there is an error term, u in the population.

So, this is the population specification of the model. But if we exclude experience square from this specification, then that violates the Gauss Markov theorem assumptions and we will not get unbiased estimators because you can see, in that case, experience square will be a part of the

error term and as a result of which the assumption that expected value of u conditional upon all sorts of values of x that is all functional forms of x equals to 0 is actually not valid anymore.

So, similarly, suppose the logarithm of wage is equal to β_0 plus β_1 education plus β_2 experience plus β_3 ability plus an error term. However, since the ability is not observed and is not included, then we have the problem of model misspecification. So, the ability is actually difficult to measure. So, it is quite possible that we can exclude this variable or it will not be included in the expression.

(Refer Slide Time: 02:39)

Problem of Omitted Variables

- This kind of misspecification is known as the problem of omitted variables.
- This is also known as underspecifying the model.
- Generalizing, if the model is
$$y = \beta_0 + \beta_1 x + \beta_2 z + u$$
then
$$E(y/x, z) = \beta_0 + \beta_1 x + \beta_2 z$$
→ not observable
- But if we run the regression only on 1 and x , then z is the omitted variable and the model to be estimated becomes
$$y = \beta_0 + \beta_1 x + (\beta_2 z + u)$$

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE 5



This kind of misspecification is known as the problem of omitted variables, this is also known as underspecifying the model. So, generalizing it, if the model is y equals β_0 plus $\beta_1 x$ plus $\beta_2 z$ where z is ability then we have expected value of y given x and z is this that is β_0 plus $\beta_1 x$ plus $\beta_2 z$ but if we run the regression only on 1 and x because z is not observable and z is actually the this is not observable. So, z is the omitted variable then the model to be estimated becomes β_0 plus $\beta_1 x$ plus $\beta_2 z$ plus u so this becomes my total error term.

(Refer Slide Time: 03:32)

Problem of Omitted Variables

- : If the composite error term ($\beta_2 z + u$) is uncorrelated to x , i.e.
 - $E(u/x) = 0$ &
 - $E(z/x) = 0$
 then there is no problem;
- But if the composite error term is correlated with x , i.e. $E(z/x) \neq 0$, then running a regression of y on only x gives,



$$E(y/x, z) = \beta_0 + \beta_1 x + \beta_2 E(z/x). \quad E(u/x) = 0.$$
- It is highly likely that ability and education are correlated such that,



6

Problem of Omitted Variables

- $E(z/x) = \theta_0 + \theta_1 x \rightarrow z = \theta_0 + \theta_1 x + u_1$ [where $\theta_1 > 0$ is expected]
- Then, $E(y/x, z) = \beta_0 + \beta_1 x + \beta_2 (\theta_0 + \theta_1 x)$

$$= (\beta_0 + \beta_2 \theta_0) + (\beta_1 + \beta_2 \theta_1) x \quad \beta_1 + \beta_2 \theta_1 \neq \beta_1$$
- Therefore, the estimated coefficient will not be β_1 ; hence, OLS estimate of β_1 , $\hat{\beta}_1$, does not converge to β_1 . If we call the estimated parameter $\tilde{\beta}_1$, such that $\tilde{\beta}_1 = \beta_1 + \beta_2 \theta_1$, then
- $E(\tilde{\beta}_1) = E(\beta_1) + E(\beta_2) \hat{\theta}_1$
- $\hat{\theta}_1$ is considered non-random because it depends only on the independent variables.



7

If the composite error term $\beta_2 z + u$ is uncorrelated to x , that is expected value of u given x is 0 and the expected value of z given x is 0, then there is absolutely no problem we still get an unbiased estimator of β_1 . But, if the composite error term is correlated with x such that the expected value of z given x is not equal to 0, that is, if ability and education are correlated with each other, then running a regression of y on x only gives this expression.

So, we will be left with β_2 expected value of z x because the expected value of z x is not equal to 0, we still assume that the expected value of u given x is 0, but since the expected value

of z is given x is not zero (*see the slide above*). So, this is actually my original expression, it is highly likely that ability and education are correlated such that the expected value of z given x equals to $\theta_0 + \theta_1 x$. So, we are assuming a linear relationship between z and x such that the expected value of z given x is equal to $\theta_0 + \theta_1 x$.

Alternatively, this can be written as Z is equal to $\theta_0 + \theta_1 x + u$ or here we can denote the error term by u_1 say, in order to distinguish it from the original model error term. So, now here θ_1 is greater than zero that is what is expected, because, we can say that as education increases ability also increases. In that case expected value of y given x and z will be, replacing the inner expression for expected values of z given x with $\theta_0 + \theta_1 x$.

So, β_2 into $\theta_0 + \theta_1 x$, expanding this is the intercept term and this is the slope term. So, in case I have omitted a variable that is correlated with the existing regressor that is x in that case, the coefficient estimate of the regressor that is existing regressor, that is x is actually not equal to β_1 . So, $\beta_1 + \beta_2 \theta_1$ and this is not equal to β_1 (*refer slide time 03:32*).



So, that is the problem that we will get in case we have an omitted variable and the variable is actually correlated with the existing regressor. Therefore, the estimated coefficient will not be β_1 . Hence, the OLS estimate of β_1 that is $\hat{\beta}_1$ does not converge to β_1 . If we call the estimated parameter $\tilde{\beta}_1$ such that $\tilde{\beta}_1$ is equal to $\beta_1 + \beta_2 \theta_1$ then the expected value of $\tilde{\beta}_1$ is equal to the expected value of β_1 plus the expected value of $\beta_2 \theta_1$.

Now, $\hat{\theta}_1$ is taken or kept outside the expected expression, because it is considered to be non-random as it depends only on the independent variables and we replace θ_1 with its estimated counterpart that is $\hat{\theta}_1$.

(Refer Slide Time: 07:11)

Problem of Omitted Variables

- Therefore, $E(\tilde{\beta}_1) = \beta_1 + \beta_2 \hat{\theta}_1$
- Therefore, bias in $\tilde{\beta}_1 = E(\tilde{\beta}_1) - \beta_1 = \beta_2 \hat{\theta}_1 \neq 0$
- $\tilde{\beta}_1$ will be unbiased if either β_2 or $\hat{\theta}_1$ or both are zero.
- Given that the original model is $y = \beta_0 + \beta_1 x + \beta_2 z + u$, if $\beta_2 = 0$, then we might not have the problem of omitted variables.
- On the other hand, even if variable z is important, so that $\beta_2 \neq 0$, but it has no correlation with x , so that $\hat{\theta}_1 = 0$, then also there will be no problem of omitted variables.

 IIT Kharagpur  NPTEL ONLINE CERTIFICATION COURSE 8

Therefore, the expected value of beta one tilde is equal to beta 1 plus beta 2 theta 1 hat therefore, the bias in beta 1 tilde is actually the expected value of beta 1 tilde minus beta 1 which is beta 2 theta 1 hat (*refer slide time: 07:11*).

So, we can see that beta one tilde will be unbiased if either beta 2 or theta 1 hat or both are 0, given that the original model is $y = \beta_0 + \beta_1 x + \beta_2 z + u$ if beta 2 is equal to 0, then we actually do not have an omitted variable problem altogether, because in that case, it shows that z is actually not explaining why. So, the problem of omitted variables is not there.

On the other hand, even if variables that are important, beta 2 is not equal to 0, but that z and x are not correlated i.e. z do not have any correlation with x . So, that theta 1 hat is equals to 0 then also there will be no problem of omitted variables. But, if a variable is actually omitted, and it is correlated with the existing regressor x then the estimated parameter beta 1 tilde is actually a biased estimator. Because in that case, our beta 2 theta 1 hat will not be equals to 0. This was the problem of the omitted variable.



(Refer Slide Time: 08:43)

Measurement Error (ME)

- ME can be there in the dependent variable as well as the independent variables.
- We would first consider ME in the **dependent variable**.
- Suppose the observed values of the dependent variable is 'y', but there is error in the measurement such that

$$e_0 = y - y^* \quad \text{or} \quad y = y^* + e_0$$
 where y^* is the actual observation and the original model is

$$y^* = \alpha + \beta x + u$$



9

Measurement Error (ME)

- But with observed values the model becomes

$$y - e_0 = \alpha + \beta x + u$$

$$\Rightarrow y = \alpha + \beta x + (u + e_0)$$
- Remember that the CLRM assumptions are



$$\text{cov}(x, e_0) = 0, \text{cov}(y^*, e_0) = 0, \text{ \& cov}(u, e_0) = 0;$$
 i.e. the ME is a purely random effect.

$$y^* = y - e_0, \text{cov}(x, y^*) = 0, \text{cov}(x, e_0) = 0, \text{cov}(u, e_0) = 0$$

$$V(u + e_0) = V(u) + V(e_0) = \sigma_u^2 + \sigma_{e_0}^2$$
- When we run a regression of y on 1 and x , the estimated coefficients remain unbiased.
- Only the standard error of $\hat{\alpha}$ and $\hat{\beta}$ goes up.

$$\sigma_u^2 (x'x)^{-1}$$

$$(\sigma_u^2 + \sigma_{e_0}^2) (x'x)^{-1}$$



10

Now, we talk about the problem of measurement error another kind of misspecification i.e. measurement error. Measurement error can be there in the dependent variable as well as in the independent variables, we would first consider measurement error in the dependent variable that is variable y . Suppose, the observed value of the dependent variable is y , but there is an error in the measurement such that $e_0 = y - y^*$ or $y = y^* + e_0$ this is the same thing written right by rearranging terms (refer slide time :08: 43; first slide).

So, the thing is that y^* is the actual observation, but we are observing the only y . So, the original model is $y^* = \alpha + \beta x + u$, but what we are observing? We are observing y . So, since $y^* = y - e$, that is why we have $y - e = \alpha + \beta x + u$ which implies $y = \alpha + \beta x + u + e$. So, now, I have a greater error term or larger error term.

Remember that the classical linear regression model assumptions are covariance between x and e should be 0. So, the covariance between x and u is anyway assumed to be 0, the covariance between x and u are assumed to be 0. Now, when we have an extended error term, so, additionally we must have covariance between x and $e = 0$.

Covariance between y^* and e is 0, the covariance between e ; and u and e is also 0 which implies that the measurement error in the dependent variable is purely a random factor or random effect. When we run a regression of y on 1 and x , the estimated coefficients remain unbiased if the measurement error is actually completely random, only the standard error of $\hat{\alpha}$ and $\hat{\beta}$ goes up (*refer slide time :08: 43; second slide*).

Now, since you know population error has gone up, population error variance is also supposed to go up even if u and e are uncorrelated with each other. The variance of $u + e$ will be a variance of u plus the variance of e covariance between u and e is 0. We can write as $\sigma_u^2 + \sigma_e^2$.

Now, if you remember that, the variance of the estimated parameters is $\sigma_u^2 + \sigma_e^2$ into $(X'X)^{-1}$. Now, in the case of this situation, I have $\sigma_u^2 + \sigma_e^2$ into $(X'X)^{-1}$. So, now, you can see that the error variance has actually gone up. (Refer Slide Time: 11:53)

Measurement Error (ME)

- The estimated parameters lose out on efficiency but still remain consistent.
- Therefore, ME in dependent variable is not a problem the long there is consistency.
- Now we consider ME in **independent variables**.
- Suppose the true model is $y = \alpha + \beta x^* + u$
- But the observed variable is $x = x^* + e_0$
- If we assume $\text{cov}(x^*, e_0) = 0$, then $V(x) = V(x^*) + V(e_0)$



The estimated parameters lose out on efficiency but still remain consistent. Therefore, measurement in the dependent variable is not a problem, as long as there is consistency. Now, we consider measurement error in the independent variables, suppose, our true model is y equals α plus βx star plus u , but what we observe is x which is x star plus e naught (*refer slide 11:53*). So, again there is a measurement error in the independent variable. Instead of observing x star, we are observing x , which is the original value plus some error term added to it. Now, there could be many situations of this kind of measurement error.

And measurement errors are many times completely unintentional. We can take the example of household savings reported by individuals. So, individuals can be different household members. The household saving reported by husbands can be different from those reported by wives and can be also reported by other family members. Similarly, we can talk about things which are reported by individuals, for example, the number of days ill in a particular year.

So, the number of days actually ill could be different from the number of days reported by an individual and that difference could be a completely random term. It is possible that I remained ill for say 10 days in a year, but instead of mentioning it at 10 days, I reported it to be 20 days. If it is intentional, then it is intentional, not only for me but for a large number of people. If we are biased towards exaggerating it, or if we have a tendency to report a lower number of days, in both cases, we would see a systematic error.



And in that case, that error is actually not completely random. But the long errors are completely random, which is also possible that I forgot actually how many days I remained ill. So, in that case, that is probably not much of a problem at least in the case, when we have a measurement error in the dependent variable.

Now, considering the case of measurement error in the independent variable, so, this is my observed model or observed variable which is x^* plus e_0 (refer slide 11:53). If we assume that covariance between x and e_0 , x^* and e_0 which is 0, then the variance of x is equal to the variance of x^* plus the variance of e_0 . Again we do not have a covariance term coming out.

(Refer Slide Time: 14:40)

Measurement Error (ME)

- CLRM assumptions require
 $\text{cov}(x^*, e_0) = 0, \text{cov}(u, e_0) = 0$ & $\text{cov}(x, e_0) = \sigma_{e_0}^2$
 since, $\text{cov}(x, e_0) = \text{cov}(x^* + e_0, e_0)$
 $= \text{cov}(x^*, e_0) + \text{cov}(e_0, e_0) = 0 + V(e_0) = \sigma_{e_0}^2$
- So, when there is ME in the independent variable then,
 $y = \alpha + \beta(x - e_0) + u = \alpha + \beta x + (u - \beta e_0)$
- Since, $\text{cov}(x, e_0) \neq 0$, we will not get consistent estimates.



12

So, the classical linear regression model assumptions required covariance of x^* and e_0 to be 0; covariance between u and e_0 to be 0 and covariance between x and e_0 is equal to $\sigma_{e_0}^2$, why this is so? Because you can see that covariance of x and e_0 is basically covariance x^* and e_0 and e_0 and e_0 . Now, I expand it, we have covariance x^* and e_0 , covariance e_0 and e_0 and covariance x^* (refer slide time 14:40).

Now covariance x star e naught is 0 and covariance e naught, e naught is actually variants of e naught, if we denote it by sigma square e naught so, that is how we have covariance x e naught is equal to sigma square e naught. So, when there is a measurement error in the independent variable, then this is our model i.e. y equals alpha plus beta x minus x naught plus u , and my error term which is u minus beta e naught (refer slide 11:53). Since covariance x e naught is not equal to 0 we will not get consistent estimates.

(Refer Slide Time: 15:53)

Measurement Error (ME)

- Proof:
- When the model is $y = \alpha + \beta x + u$
- We know that $\hat{\beta} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}$
- Also we know that $(y_i - \bar{y}) = \alpha + \beta x_i - \alpha - \beta \bar{x} + u_i - \bar{u} = \beta(x_i - \bar{x}) + u_i - \bar{u}$
- Substituting $(y_i - \bar{y})$ into $\hat{\beta}$, and rearranging terms we obtain
- $\hat{\beta} = \frac{\sum \beta(x_i - \bar{x})^2 + \sum(x_i - \bar{x})(u_i - \bar{u})}{\sum(x_i - \bar{x})^2} = \beta + \frac{\sum(x_i - \bar{x})(u_i - \bar{u})/n}{\sum(x_i - \bar{x})^2/n} = \beta + \frac{\text{Cov}(x, u)}{\text{Var}(x)}$

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE 13

Measurement Error (ME)

- Therefore, $\hat{\beta} \xrightarrow{P} \beta + \frac{\text{Cov}(x, u)}{\text{Var}(x)} = 0$ since $\text{Cov}(x, u) = 0$
- Now when there is ME in the independent variable then,
- $\hat{\beta} \xrightarrow{P} \beta + \frac{\text{Cov}(x, u - \beta e_0)}{\text{Var}(x)}$ since $\text{Cov}(x, u) = 0$ and $x = x^* + e_0$
- We obtained that $\text{Cov}(x, e_0) = \sigma_{e_0}^2$
- Hence, $\hat{\beta} \xrightarrow{P} \beta + \frac{-\beta \sigma_{e_0}^2}{\sigma_{x^*}^2 + \sigma_{e_0}^2}$
- Or $\hat{\beta} \xrightarrow{P} \frac{\beta \sigma_{x^*}^2}{\sigma_{x^*}^2 + \sigma_{e_0}^2}$

IIT ROORKEE NPTEL ONLINE CERTIFICATION COURSE 14

So, this is proof now, that when the model is y equals α plus βx plus u , that is the standard model, then we know that this is the expression of our $\hat{\beta}$, which is basically if I write it in matrix format, it would be $x'x^{-1}x'y$. Otherwise, this is summation $x_i - \bar{x}$ into $y_i - \bar{y}$ divided by summation $x_i - \bar{x}$ whole square (*refer to slide time 15:53; first slide*). This has been already derived and mentioned several times so far.

Also, we know that $y_i - \bar{y}$ is equal to $\alpha + \beta x_i + u_i - \alpha - \beta \bar{x} - u$ corresponds to this y . Now, by connecting terms, we can see that $\alpha - \alpha$ cancels out, I collect β . So, β has $\beta x_i - \beta \bar{x} + u_i - u$ (*refer slide 15:53; first slide*). Now, if I substitute this expression into the expression for the $\hat{\beta}$ and rearrange terms, then you can see that, first of all, I will be having β into $x_i - \bar{x}$ whole square that is the expression is put here plus $x_i - \bar{x}$ into $y_i - \bar{y}$ (*refer slide 15:53; first slide*).

The denominator remains as it is, these two things cancel out (*refer slide 15:53; first slide*). I have β here plus the entire thing divided by the denominator. Now, I am dividing both the numerator and the denominator by n . And what I am getting? In the numerator, I have covariance xu and in the denominator, I have a variance of x . Therefore, in the probability limit, $\hat{\beta}$ in probability limit tends to $\beta + \text{covariance } xu \text{ divided by variance of } x$.

Now, when there is a measurement error in the independent variable, then we have $\hat{\beta}$ with probability limit tending to $\beta + \text{covariance between } x \text{ and } u - \beta e$ because now, my error is not u , it is $u - \beta e$. You can see that the covariance between x and u equals to 0 this term is equal to 0 and the $\hat{\beta}$ in probability limit tends to β . So, we say that the estimates are consistent (*refer slide 15:53; second slide*).

Now, we need to see whether this term tends to 0 or not or this term is equal to 0 or not. So, we know that covariance between x and e is σ_e^2 , e square. So, hence $\hat{\beta}$ in probability limit tends to $\beta + \text{covariance between } x \text{ and } u - \beta e$. We can expand this, that is the covariance between $xu - \beta e$ as the covariance between x



and u minus beta covariance between x and e naught right this is 0 and this is equal to sigma square e naught (refer slide 15:53; second slide).

Therefore, we have in the numerator minus beta sigma u naught square divided by variance of x was observed to be sigma x squared plus sigma e naught squared. Therefore, probability beta hat with probability limit tends to rearranged term and it is not actually tending to beta.

(Refer Slide Time: 19:34)

Measurement Error (ME)

- The estimator is not consistent. It is called attenuation bias.
- Since $\frac{\sigma_{x^*}^2}{\sigma_{x^*}^2 + \sigma_{e_0}^2} < 1$ $\hat{\beta} < \frac{\beta \sigma_{x^*}^2}{\sigma_{x^*}^2 + \sigma_{e_0}^2}$
- Therefore, if estimated coefficients are very small, we may suspect presence of ME in x . If ME is massive then $\hat{\beta} \rightarrow 0$.
- In this context, it is important to note that
- If $Cov(x_i, u_i) \neq 0$, then x_i is endogenous regressor.
- If $Cov(x_i, u_i) = 0$, then x_i is contemporaneously exogenous.
- If $Cov(x_i, u_j) = 0$, then x_i is strict exogenous. $Cov(x_i, u_j) \neq 0$



15

The estimator is not consistent and it is called attenuation bias. Since this expression is less than 1 because you can see all of them are positive numbers and the numerator is smaller than the denominator. So, we have this expression less than 1 which implies that beta hat will be less than this (refer slide time: 19:34). Therefore, if estimated coefficients are very small, the beta hat is

actually less than beta multiplied by a fraction. So, if estimated coefficients are very small, we may suspect the presence of measurement error in x.

If the measurement error is massive, then the beta hat pretends to 0 because this number actually becomes very small. In this context, it is important to note that if covariance between x_i and u_i is not equal to 0, then x_i is endogenous regressor or we call x_i to be an endogenous regressor. Endogenous is something that is explained by the system. So, generally we assume independent variables to be non-random, or they are independent in the sense that the model is not trying to explain them.

But if they are correlated with the error term, then this implies that the model contains additional information that may explain the variations in x as well. That is why we call them endogenous regressors. If x_i and u_i are equal to 0 or covariance between x_i and u_i equals to 0, then x_i is said to be contemporaneously exogenous that is this is a situation where covariance between x_i and u_j not equal to 0, only this is equals to 0.

So, in that case we say that x_i is contemporaneously exogenous and if covariance x_i and u_j are 0. That is all the x_i 's and all the u_i 's for any i and j is 0, then x is strict exogenous that is the model does not explain or contain any information that may explain the independent variable x.

(Refer Slide Time: 21:48)

Measurement Error (ME) in Non-Classical Set up

- Suppose h_i^* : number of days ill
 x_i : income of i
 h_i : reported number of days ill
- Model: $h_i^* = \alpha + \beta x_i + u_i$ $\beta < 0$
- Let $h_i = h_i^* + e_i$ & $\text{cov}(e_i, x_i) > 0$ *income ↑ error ↑*
 such that $E(e/x) = \theta_0 + \theta_1 x$, $\theta_1 > 0$
- There is violation of classical assumption. Now, the observed model is
 $h_i = \alpha + \beta x_i + (u_i + e_i)$
 $E(h/x) = \alpha + \beta x + E[(u_i + e_i)/x]$

Measurement Error (ME) in Non-Classical Set up

- Substituting for $E[(u + e)/x]$, we obtain
- $E(h/x) = \alpha + \beta x + \theta_0 + \theta_1 x = (\alpha + \theta_0) + (\beta + \theta_1)x$ $\beta < 0$
- Since $\beta < 0$ and $\theta_1 > 0$, $(\beta + \theta_1)$ can be positive or negative.
- But if ME is very large then $(\beta + \theta_1) > 0$ and the regression will not report expected results. $|\beta| < |\theta_1|$
- We will have biased parameter estimates.
- Therefore, the estimates are neither best nor consistent.

Now, we talked about measurement errors in the non-classical setup. Suppose h_i^* is the number of days ill, I will explain why we call it a non-classical setup. Now, x_i is the income of i . And h_i is the reported number of days an individual remains ill (*refer slide 21:48*).

So, the model is actually the number of days ill is explained as a function of the income of individual i . When we expect the beta to be less than 0, that is, as income increases, the number of days a person remain ill actually reduces, but this is actually the reported number of days ill which is different from or which could be different from the actual number of days the person remains ill.

And suppose this is how they are related. So, h_i is equals to $h_i^* + e_i$. And now, we are assuming that the covariance between this error term and income, are correlated, which implies, if this is positive, then this implies that as income increases the tendency or the error actually increases, and as income increases, the error also increases such that expected value of E given x is $\theta_0 + \theta_1 x$ where θ_1 is greater than 0 (*refer slide 21:48; first slide*).

So, we are expecting a positive correlation and that is why a positive parameter estimate. If we have a linear relationship between e and x , the error here and the independent variable there is a violation of classical assumption. So, that is why this is called non classical setup, because now, the observed model is $h_i = \alpha + \beta x_i + u_i + e_i$, where this x_i and u_i are correlated.

So, the expected value of h given x equals α plus βx plus the expected value of u plus e given x .

Substituting expected value of u plus e given x we obtain the expected value of h given x equals α plus βx plus θ_0 plus $\theta_1 x$. So, this expression, this is the linear relation we obtain between e and x and that is what is being replaced here as a result of which this is my constant term and this is my parameter of x (refer slide 21:48; second slide).


Now, this parameter of x is actually not equal to β . Since β is expected to be less than 0 and θ_1 is greater than 0. $\beta + \theta_1$ can be either positive or negative. But if the measurement error is very large, then we expect θ_1 to be large. And as a result of which mod value of β could be less than the mod value of θ_1 and as a result of which $\beta + \theta_1$ will be greater than 0.

And the regression will not report expected results because originally we expected β to be less than 0 if we observe a coefficient of x which is greater than 0, then this is actually not fulfilling our expectations following the theoretical argument or for what we have initially postulated. So, we will have biased parameter estimates. Therefore, the estimates are neither best nor consistent. So they have neither minimum variance nor consistency.

(Refer Slide Time: 25:31)

RESET

- RESET is regression specific error test suggested by Ramsey in 1969 as general test for functional misspecification.
- The idea behind RESET is if the original model is
$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + u \quad (1)$$
- And it satisfies $E(u|x_1, \dots, x_k) = 0$, then no non-linear functions of the independent variables should be significant when added to the equation (1).
- RESET adds polynomials in the OLS fitted values to equation (1), and most often the squared and cubed terms of the fitted values have proven to be useful.
- Let \hat{y} denotes the OLS fitted values from estimating equation (1).

IIT KHARAGPUR
NPTEL ONLINE
CERTIFICATION COURSE18

RESET

: The expanded equation becomes

$$\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \delta_1 \hat{y}^2 + \delta_2 \hat{y}^3 + u \quad (2)$$

- Equation (2) tests whether equation (1) has missed important nonlinearities or not.
- The null hypothesis is equation (1) is correctly specified; i.e.
- $H_0: \delta_1 = 0, \delta_2 = 0$
- If the null hypothesis is rejected, it suggests some sort of functional form problem.

Now, how do we test whether there is misspecification or not? RESET is regression specific error test. RESET was suggested by Ramsey in 1969 as a general test for functional misspecification. The idea behind reset is that, if the original model is $y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + u$, I consider a general model with k plus 1 variable and it satisfies the expected value of u given x_1 to x_k equals to 0 that is there is no correlation between the error term and the independent variables, then no nonlinear functions of the independent variables should be significant when added to the equation 1 (*refer slide time 25:31; first slide*).

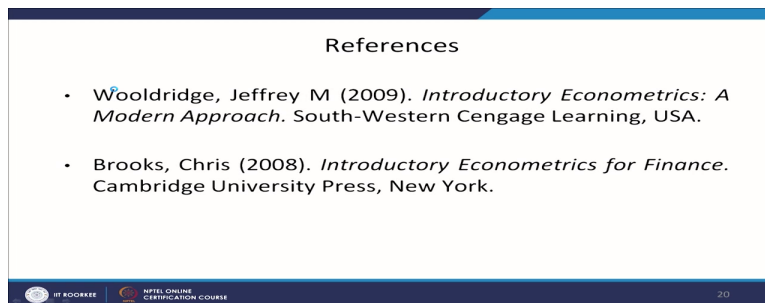
So, now, we are basically talking about or trying to find out a model misspecification where the model could have some nonlinear functions of the independent variables. That is an example, we gave in the beginning that there could be experienced quite in the specification also. So, whether that experience square has been missed out or not could be given by reset. Reset adds polynomials in the OLS fitted values to equation 1, this is our equation 1 (*refer slide time 25:31; first slide*).

And most often the squared and cubed terms of the fitted values have proven to be useful. So, that is why instead of considering a large number of polynomials, it at the max considers the squared and cubed terms of the independent variables or the fitted values. Let \hat{y} denotes OLS fitted values from estimates obtained from equation 1, the expanded equation then becomes the

original equation plus squared term of the fitted values plus a cubic term of the fitted values plus the error term.

So, equation 2 tests whether equation 1 has missed important nonlinearities or not, the null hypothesis in equation 1 is correctly specified as $\delta_1 = 0$ and $\delta_2 = 0$ (*refer slide time 25:31; second slide*). So, we are actually focusing on these two parameter estimates. If these two parameter estimates are equal to 0, then there is actually no model misspecification in terms of exclusion of nonlinear functional forms of the independent variables, existing independent variables. If the null hypothesis is rejected, it suggests some sort of functional form problem i.e. some sort of functional or nonlinearities in one or more independent variables would have been present in the model specification. But if the null hypothesis is not rejected, then probably we do not have any problem.

(Refer Slide Time: 28:21)



So, these are the references. Thank you.