

Homework 3

Tate Mason

An ECON - 8070 Homework Assignment

November 13, 2024

Question 4.2

Problem

(a)

Show that, under random sampling and the zero conditional mean assumption $\mathbb{E}(u|\mathbf{x}) = 0$, $\mathbb{E}(\hat{\beta}|\mathbf{X}) = \beta$ if $\mathbf{X}\mathbf{X}'$ is nonsingular. (Hint: use property CE.5 in the appendix of chapter 2)

(b)

In addition to the assumption from part a, assume that $\text{var}(u|\mathbf{x}) = \sigma^2$. Show that $\text{var}(\hat{\beta}|\mathbf{X}) = \sigma^2((\mathbf{X}\mathbf{X}')^{-1})$

Solutions

(a)

Proof. The OLS estimator is given by $\hat{\beta} = (X'X)^{-1}X'y$ and the linear model gives us that $y = X\beta + u$. If we substitute that into the formula for $\hat{\beta}$, we get $\hat{\beta} = (X'X)^{-1}X'X\beta + (X'X)^{-1}X'u$. Next, taking the conditional expectation of $\hat{\beta}$ given X and using the zero mean assumption, we can show that $\mathbb{E}(\hat{\beta}|X) = (X'X)^{-1}X\beta + (X'X)^{-1}X'0 \rightarrow \mathbb{E}(\hat{\beta}|X) = (X'X)^{-1}X'X\beta \rightarrow \mathbb{E}(\hat{\beta}|X) = \beta$ \square

(b)

Proof. The formula for the variance of the OLS estimator $\hat{\beta} = \mathbb{E}[(\hat{\beta} - \mathbb{E}(\hat{\beta}|X))(\hat{\beta} - \mathbb{E}(\hat{\beta}|X))'|X']$. Using the result from (a), we can say that $\text{var}(\hat{\beta}|X) = \mathbb{E}[((X'X)^{-1}X'u)((X'X)^{-1}X'u)'|X] = (X'X)^{-1}X'\mathbb{E}(uu'|X)X(X'X)^{-1}$. As the problem tells us that $\text{var}(u|X) = \sigma^2$, we can state $\mathbb{E}(uu'|X) = \sigma^2 I$. Thus, $\text{var}(\hat{\beta}|X) = (X'X)^{-1}X'(\sigma^2 I)X(X'X)^{-1} = \sigma^2(X'X)^{-1}$ \square

question 4.3

problem

suppose that in the linear model 4.5, $(\mathbf{x}'u) = 0$ (where \mathbf{x} contains unity), $\text{var}(u|\mathbf{x}) = \sigma^2$, but $(u|\mathbf{x}) \neq (u)$.

(a)

is it true that $(u^2|\mathbf{x}) = \sigma^2$?

(b)

what relevance does part a have for ols estimation?

solutions

(a)

no, this is not true.

Proof. $(\text{var}u|x) = (u^2|x) - (u|x)^2$. from the question, we can say that $\sigma^2 = (u^2|x) - (u|x)^2$. due to the fact that $(u|x) \neq (u)$, we know that u varies with x , and the same with $(u|x)^2$. so, rearranging the equation, $(u^2|x) = \sigma^2 + (u|x)^2$. since, as stated, $(u|x)^2$ varies with x , so too must $(u^2|x)$. with that conclusion we can say that $(u^2|x) \neq \sigma^2$ \square

(b)

when heteroskedasticity is present, the usual ols s.e. will be incorrect due to the assumption of homoskedasticity. thus, robust s.e. will be needed to produce meaningful estimation with ols.

question 4.17

problem

consider the standard linear model $y = \mathbf{x}\boldsymbol{\beta} + \mathbf{u}$ under assumptions ols.1 and ols.2. define $h(\mathbf{x}) \equiv [u^2|\mathbf{x}]$. let $\hat{\boldsymbol{\beta}}$ be the ols estimator, and show that we can always write:

$$\text{avar}\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = [[\mathbf{x}'\mathbf{x}]]^{-1} [h(\mathbf{x})\mathbf{x}'\mathbf{x}] [[\mathbf{x}'\mathbf{x}]]^{-1} \quad (1)$$

this expression is useful when $[u|\mathbf{x}] = \mathbf{0}$ for comparing the asymptotic variances of ols and weighted least squares estimators; see, for example, wooldridge (1994b).

solution

Proof. the asymptotic distribution of the ols distribution is given as $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = (n^{-1} \sum_{i=1}^n \mathbf{x}'_i \mathbf{x}_i)^{-1} (n^{-\frac{1}{2}} \sum_{i=1}^n \mathbf{x}'_i u_i)$. by central limit theorem, $n^{-\frac{1}{2}} \sum_{i=1}^n \mathbf{x}'_i u_i \xrightarrow{d} N(0, [(u^2|x)\mathbf{x}'\mathbf{x}])$ which as given in the hint, $(u^2|x) \equiv h(\mathbf{x})$. by law of large numbers, $n^{-1} \sum_{i=1}^n \mathbf{x}'_i \mathbf{x}_i \xrightarrow{p} [\mathbf{x}'\mathbf{x}]$. after plugging this into common knowledge, we can say that $\text{avar}\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = (\mathbf{x}'\mathbf{x})^{-1} (h(\mathbf{x})\mathbf{x}'\mathbf{x}) (\mathbf{x}'\mathbf{x})^{-1}$. \square

question 5.1

problem

in this problem you are to establish the algebraic equivalence between 2sls and ols estimation of an equation containing an additional regressor. although the result is completely general, for simplicity consider a model with a single (suspected) endogenous variable:

$$\begin{aligned} y_1 &= z_1\delta_1 + \alpha_1 y_2 + u_1, \\ y_2 &= z_2\pi_2 + v_2. \end{aligned}$$

for notational clarity, we use y_2 as the suspected endogenous variable and z as the vector of all exogenous variables. the second equation is the reduced form for y_2 . assume that z has at least one more element than z_1 . we know that one estimator of (δ_1, α_1) is the 2sls estimator using instruments x . consider an alternative estimator of (δ_1, α_1) : (a) estimate the reduced form by ols, and save the residuals \hat{v}_2 ; (b) estimate the following equation by ols:

$$y_1 = z_1\delta_1 + \alpha_1 y_2 + \rho \hat{v}_2 + \text{error}. \quad (5.52)$$

show that the ols estimates of δ_1 and α_1 from this regression are identical to the 2sls estimators. (hint: use the partitioned regression algebra of ols. in particular, if $\hat{y} = x_1\beta_1 + x_2\beta_2$ is an ols regression, β_1 can be obtained by first regressing x_1 on x_2 , getting the residuals, say \hat{x}_1 , and then regressing y on \hat{x}_1 ; see, for example, davidson and mackinnon (1993, section 1.4). you must also use the fact that z_1 and \hat{v}_2 are orthogonal in the sample.)

solution

Proof. 2sls in this case would work by first regressing y_2 on z to get \hat{y}_2 . then, we would regress y_1 on z_1 and \hat{y}_2 . for the alternative method, $\hat{v}_2 = y_2 - z\hat{\pi}_2$ which gives us the residuals from the reduced form problem. then, we would run ols on equation 5.52. to derive δ_1 and α_1 , we regress y_2 and z_1 on \hat{v}_2 and get the residuals. then, we regress y_1 on the residuals from the last step. since, as given, \hat{v}_2 is orthogonal to z_1 , the residual from the regression of z_1 on \hat{v}_2 is z_1 . the residual from the regression of y_2 on \hat{v}_2 is \hat{y}_2 . this is the same result as when regressing y_1 on z_1 and \hat{y}_2 . the two estimates, therefore, show identical estimates for δ_1 and α_1 . \square

question 5.2

problem

consider a model for the health of an individual:

$$\begin{aligned} \text{health} = & \beta_0 + \beta_1 \text{age} + \beta_2 \text{weight} + \beta_3 \text{height} \\ & + \beta_4 \text{male} + \beta_5 \text{work} + \beta_6 \text{exercise} + u_1, \end{aligned} \quad (2)$$

where *health* is some quantitative measure of the person's health; *age*, *weight*, *height*, and *male* are self-explanatory; *work* is weekly hours worked; and *exercise* is the hours of exercise per week.

parts

(a)

why might you be concerned about *exercise* being correlated with the error term u_1 ?

(b)

suppose you can collect data on two additional variables, *disthome* and *distwork*, the distances from home and from work to the nearest health club or gym. discuss whether these are likely to be uncorrelated with u_1 .

(c)

now assume that *disthome* and *distwork* are in fact uncorrelated with u_1 , as are all variables in equation (??) with the exception of *exercise*. write down the reduced form for *exercise*, and state the conditions under which the parameters of equation (??) are identified.

(d)

how can the identification assumption in part c be tested?

solutions

(a)

things like health conditions can affect ability to exercise as well as those who care more about health will be more likely to exercise and be healthier. these unobserved factors, as well as the potential for bad health to effect ability exercise, are reasons why exercise may be correlated with the error term.

(b)

these are likely uncorrelated factors as the choice of a dwelling or a workplace is typically decided by more factors than just proximity to a gym. thus, it would be a valid addition to the model.

(c)

the reduced form equation is as follows:

$$exercise = \pi_0 + \pi_1 age + \pi_2 weight + \pi_3 height + \pi_4 male + \pi_5 work + \pi_6 disthome + \pi_7 distwork + v_2$$

identification has three primary conditions:

a) rank condition states that one of π_6, π_7 must be non-zero

b) the exclusion restriction is assumed in this question

c) overidentified equation as the instruments outnumber endogenous variables

(d)

one example of a test for the identification assumption in (c) is the first stage f-test. tests joint significance of instruments in reduced form.

question 5.11

problem

a model with a single endogenous explanatory variable can be written as

$$y_1 = z_1\delta_1 + \alpha_1 y_2 + u_1, \quad (3)$$

$$e(z'u_1) = 0, \quad (4)$$

where $z = (z_1, z_2)$. consider the following two-step method, intended to mimic 2sls:

(a)

regress y_2 on z_2 , and obtain fitted values, \hat{y}_2 . (that is, z_1 is omitted from the first-stage regression.)

(b)

regress y_1 on z_1, \hat{y}_2 to obtain $\hat{\delta}_1$ and $\hat{\alpha}_1$. show that $\hat{\delta}_1$ and $\hat{\alpha}_1$ are generally inconsistent. when would $\hat{\delta}_1$ and $\hat{\alpha}_1$ be consistent? (hint: let y_2^0 be the population linear projection of y_2 on z_2 , and let a_2 be the projection error: $y_2^0 = z_2\lambda_2 + a_2$, $e(z'a_2) = 0$. for simplicity, pretend that λ_2 is known rather than estimated; that is, assume that \hat{y}_2 is actually y_2^0 . then, write)

$$y_1 = z_1\delta_1 + \alpha_1 y_2^0 + \alpha_1 a_2 + u_1 \quad (5)$$

and check whether the composite error $\alpha_1 a_2 + u_1$ is uncorrelated with the explanatory variables.

solution

Proof. as given in the hint, $y_2^0 = z_2\lambda_2 + a_2$ such that $\mathbb{E}(z_2'a_2) = 0$. substituting this into the main equation, $y_1 = z_1\delta_1 + \alpha_1(z_2\lambda_2 + a_2) + u_1 = z_1\delta_1 + \alpha_1 z_2\lambda_2 + (\alpha_1 a_2 + u_1)$. for consistency, $\alpha_1 a_2 + u_1$ needs to be uncorrelated with both z_1 and y_2^0 . from assumptions given in the question, we know that $(z_1'u_1) = 0$. however, $(z_1'a_2) \neq 0$ due to the exclusion of z_1 from the linear projection. with this in my mind, we can say that $\hat{\delta}_1, \hat{\alpha}_1$ are generally inconsistent. they would be consistent if z_1, z_2 are uncorrelated, y_2 is unrelated to z_1 , or if $\alpha_1 = 0$. important takeaways are that due to the omission of z_1 from the first stage is akin to incorrectly excluding relevant instruments, creating correlation between a_2 and the exogenous variables. thus, it is necessary to use valid instruments in the first stage. \square