ECON 251

Discussion Section

Week 7 Solutions

 \Rightarrow

1. Prove that OLS is unbiased under independence, and consistent under exogeneity.

$$\begin{split} \hat{\beta}_{1}^{OLS} &\coloneqq \frac{\widehat{\operatorname{Cov}}(X,Y)}{\widehat{\operatorname{Var}}(X)} = \frac{\sum_{i=1}^{N} (X_{i} - \overline{X})(Y_{i} - \overline{Y})}{\sum_{i=1}^{N} (X_{i} - \overline{X})^{2}} = \frac{\sum_{i=1}^{N} (X_{i} - \overline{X}) \cdot Y_{i}}{\sum_{i=1}^{N} (X_{i} - \overline{X})^{2}} \\ &= \frac{\sum_{i=1}^{N} (X_{i} - \overline{X}) \cdot (\beta_{0} + \beta_{1}X_{i} + U_{i})}{\sum_{i=1}^{N} (X_{i} - \overline{X})^{2}} \\ &= \frac{\beta_{0} \sum_{i=1}^{N} (X_{i} - \overline{X}) + \beta_{1} \sum_{i=1}^{N} (X_{i} - \overline{X}) X_{i} + \sum_{i=1}^{N} (X_{i} - \overline{X}) U_{i}}{\sum_{i=1}^{N} (X_{i} - \overline{X})^{2}} \\ &= \frac{\beta_{0} \cdot 0 + \beta_{1} \sum_{i=1}^{N} (X_{i} - \overline{X})^{2} + \sum_{i=1}^{N} (X_{i} - \overline{X}) U_{i}}{\sum_{i=1}^{N} (X_{i} - \overline{X})^{2}} \\ &= \beta_{1} + \frac{\sum_{i=1}^{N} (X_{i} - \overline{X}) U_{i}}{\sum_{i=1}^{N} (X_{i} - \overline{X})^{2}} \\ E[\hat{\beta}_{1}^{OLS}|X_{i}] &= E\left[\beta_{1} + \frac{\sum_{i=1}^{N} (X_{i} - \overline{X}) U_{i}}{\sum_{i=1}^{N} (X_{i} - \overline{X})^{2}} X_{i}\right] \\ &= \beta_{1} + \frac{\sum_{i=1}^{N} (X_{i} - \overline{X}) E[U_{i}|X_{i}]}{\sum_{i=1}^{N} (X_{i} - \overline{X})^{2}} = \beta_{1} \end{split}$$

$$\begin{aligned} \min_{N \to \infty} \hat{\beta}_1^{OLS} &\coloneqq \min_{N \to \infty} \frac{\widehat{\operatorname{Cov}}(X, Y)}{\widehat{\operatorname{Var}}(X)} = \frac{\min_{N \to \infty} \operatorname{Cov}(X, Y)}{\min_{N \to \infty} \widehat{\operatorname{Var}}(X)} = \frac{\operatorname{Cov}(X, Y)}{\operatorname{Var}(X)} \\ &= \frac{\operatorname{Cov}(X, \beta_0 + \beta_1 X + U)}{\operatorname{Var}(X)} = \frac{\operatorname{Cov}(X, \beta_0) + \operatorname{Cov}(X, \beta_1 X) + \operatorname{Cov}(X, U)}{\operatorname{Var}(X)} \\ &= \frac{0 + \beta_1 \cdot \operatorname{Var}(X, X) + \operatorname{Cov}(X, U)}{\operatorname{Var}(X)} = \beta_1 + \frac{\operatorname{Cov}(X, U)}{\operatorname{Var}(X)} = \beta_1 \end{aligned}$$

 Relate the statistical concepts of unbiasedness and consistency with the econometric concepts of selection bias (under the potential outcomes/treatment effects framework) and omitted variable bias (under the linear regression framework).

In the previous question, we showed that the independence assumption E[U|X] = 0implies OLS is unbiased $E[\hat{\beta}_1^{OLS}] = \beta_1$ while the exogeneity assumption Cov(X, U) = 0implies OLS is consistent $\lim_{N\to\infty} \hat{\beta}_1^{OLS} = \beta_1$. In the potential outcomes framework, the independence assumption means that we can identify the average treatment effect $ATE := E[Y_i(1) - Y_i(0)]$ with simple comparisons¹ as there is no selection bias on average between treated and control units. Thus, we write

$$\overline{Y}_1 - \overline{Y}_0 = ATT + SB = ATE$$

where the first equality is <u>always</u> true (mean differences equal causal effect of treatment among the treated plus selection bias $SB \coloneqq E[Y_i(0)|X_i = 1] - E[Y_i(0)|X_i = 0]$) and the second equality is <u>only</u> true under independence $X_i \perp Y_i(0), Y_i(1)$ which is equivalent to E[U|X] = 0. In the linear regression framework, we can imagine what might happen if we omit an important (for determining the outcome) and relevant (for correlating with treatment) variable from our model. The probability limit of OLS is always given by

$$\lim_{N \to \infty} \hat{\beta}_1^{OLS} = \beta_1 + \frac{\operatorname{Cov}(X, U)}{\operatorname{Var}(X)} = \beta_1 + \gamma^{X \to V} \cdot \delta^{V \to Y}$$

where the second equality follows from omitted variable bias (OVB) formula, with $\gamma^{X \to V}$ giving the linear effect of the treatment on the omitted variable V and $\delta^{V \to Y}$ giving its effect on the outcome. The exogeneity assumption Cov(X, U) = 0 rules out such omitted random variables which are both important and relevant, and delivers consistency!

¹ Recall that OLS is equivalent to a simple comparison (mean difference) when treatment is binary!

- 3. In *Stata*, review interpretation of simple and multiple linear regression tables. Discuss functional form issues (and possibilities) with linear regression: logarithm of outcome or treatment, binary variables and groups, quadratic forms, interactions.
- 4. Discuss "On Binscatter" by Cattaneo et al. (2022) as a rigorous method to visually inspect assumptions SLR1/MLR1, that linear regressions are actually linear!

The linearity assumption requires that the conditional expectation function relating the outcome to the treatment be linear $E[Y|X = x] = g(x) = \beta_0 + \beta_1 x$. Sometimes this assumption is overlooked² in practice, but an easy check is to simply plot the data! If the scatterplot displays non-linear patterns, then we could either consider some transformations to make things linear (such as taking logarithms of skewed variables in our analysis or adding the square of treatment to our specification). A quite practical limitation of this simple and intuitive check is when the sample size is large! In these situations, it can be quite difficult to discern the shape of the conditional expectation function, so a simple solution involves making a scatter plot of the treatment vs the average value of the outcome within certain intervals of the treatment distribution. Cattaneo et al. (2022) present methods on implementing this in practice, which involves several important technical steps like figuring out the right way to specify the bins across treatment to average the outcome and creating confidence intervals properly.

² This is not without justification: we can show that even if g(x) is non-linear, estimating slopes and an intercept with least squares produces the best linear predictor of the conditional expectation function.