

# QF Problem Set 5

a. consider the following population linear regression models

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + u$$

where  $E(u) = 0$  and  $\text{cov}(X_l, u) = 0$  for  $l \in \{1, \dots, k\}$

$$Y = \gamma_0 + \gamma_1 X_1 + \dots + \gamma_{k-1} X_{k-1} + v$$

where  $E(v) = 0$  and  $\text{cov}(X_l, v) = 0$  for  $l \in \{1, \dots, k-1\}$

$$X_k = \pi_0 + \pi_1 X_1 + \dots + \pi_{k-1} X_{k-1} + \tilde{X}_k$$

where  $E(\tilde{X}_k) = 0$  and  $\text{cov}(X_l, \tilde{X}_k) = 0$  for  $l \in \{1, \dots, k-1\}$

$$X_1 = \phi_0 + \phi_2 X_2 + \dots + \phi_{k-1} X_{k-1} + \tilde{X}_1$$

where  $E(\tilde{X}_1) = 0$  and  $\text{cov}(X_l, \tilde{X}_1) = 0$  for  $l \in \{2, \dots, k-1\}$

by Full theorem,

$$\gamma_1 = \text{cov}(Y, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

$$= \text{cov}(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + u, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

$$= \text{cov}(\beta_1 X_1 + \dots + \beta_k X_k, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

since  $\beta_0$  is a constant

and  $\tilde{X}_1$  is a linear function of  $X_1, \dots, X_k$ ,

and  $\text{cov}(X_l, u) = 0$  for  $l \in \{1, \dots, k\}$

$$= \text{cov}(\beta_1 X_1, \tilde{X}_1) / \text{var}(\tilde{X}_1) + \text{cov}(\beta_k X_k, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

since  $\text{cov}(X_l, \tilde{X}_1) = 0$  for  $l \in \{2, \dots, k-1\}$

$$= [\beta_1 \text{cov}(X_1, \tilde{X}_1) + \beta_k \text{cov}(X_k, \tilde{X}_1)] / \text{var}(\tilde{X}_1)$$

$$= [\beta_1 \text{var}(\tilde{X}_1) + \beta_k \text{cov}(X_k, \tilde{X}_1)] / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

since  $\text{cov}(X_l, \tilde{X}_1) = 0$  for  $l \in \{2, \dots, k-1\}$

$$= \beta_1 + \beta_k \text{cov}(X_k, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

$$= \beta_1 + \beta_k \pi_1$$

by Full theorem

$$b. \pi_1 = \text{cov}(X_k, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

$$= \text{cov}(X_k, X_1 - \phi_0 - \phi_2 X_2 - \dots - \phi_{k-1} X_{k-1}) / \text{var}(\tilde{X}_1)$$

$$= [\text{cov}(X_k, X_1) - \phi_2 \text{cov}(X_k, X_2) - \dots - \phi_{k-1} \text{cov}(X_k, X_{k-1})] / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator,

since  $\phi_0$  is a constant

$$= \text{cov}(X_k, X_1) / \text{var}(\tilde{X}_1)$$

given that  $\text{cov}(X_k, X_l) = 0$  for  $l \in \{2, \dots, k-1\}$

$> 0$

given that  $\text{cov}(X_k, X_1) > 0$

since  $\text{var}(\tilde{X}_1) > 0$

$$\gamma_1 = \beta_1 + \beta_k \pi_1 > \beta_1$$

given that  $\beta_k > 0$

since  $\pi_1 > 0$



2. consider the population linear regression model. I struggled quite a bit with q2,3, could you help check if I got these right?

$$Y = \gamma_0 + \gamma_1 X_1 + \gamma_2 X_2 + v$$

$$\text{where } E(v) = 0 \text{ and } \text{cov}(X_1, v) = \text{cov}(X_2, v) = 0$$

$$X_1 = \beta_0 + \beta_2 X_2 + \tilde{X}_1$$

$$\text{where } E(\tilde{X}_1) = 0 \text{ and } \text{cov}(X_2, \tilde{X}_1) = 0$$

a. Full theorem

$$\gamma_1 = \text{cov}(Y, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

$$= \text{cov}(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + u, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

$$= \text{cov}(\beta_1 X_1 + \beta_2 X_2 + u, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

since  $\beta_0$  is a constant

$\tilde{X}_1$  is a linear function of  $X_1$  and  $X_2$ ,

and  $\text{cov}(X_i, \tilde{X}_1) = 0$  for  $i \in \{1, 2\}$

$$= \text{cov}(\beta_1 X_1 + u, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

since  $\beta_0$  is a constant

and  $\text{cov}(X_2, \tilde{X}_1) = 0$

$$= \text{cov}(\beta_1 X_1 + u, X_1 - \beta_0 - \beta_2 X_2)$$

$$= \text{cov}(\beta_1 X_1 + u, X_1 - \beta_0 - \beta_2 X_2) / \text{var}(\tilde{X}_1)$$

$$= \text{cov}(\beta_1 X_1 + u, X_1 - \beta_2 X_2) / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

since  $\beta_0$  is a constant

$$= [\beta_1 \text{var}(X_1) - \beta_1 \beta_2 \text{cov}(X_1, X_2) - \beta_2 \text{cov}(u, X_2)] / \text{var}(\tilde{X}_1)$$

$$= \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

$$= \text{cov}(\beta_1 (X_1 - \beta_0 - \beta_2 X_2) + \beta_1 \beta_0 + \beta_1 \beta_2 X_2 + u, \tilde{X}_1) / \text{var}(\tilde{X}_1)$$

$$= [\beta_1 \text{var}(\tilde{X}_1) + \text{cov}(\beta_1 \beta_0 + \beta_1 \beta_2 X_2 + u, \tilde{X}_1)] / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

$$= \beta_1 + \text{cov}(\beta_1 \beta_2 X_2 + u, X_1 - \beta_0 - \beta_2 X_2) / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

since  $\beta_1 \beta_0$  is constant

$$= \beta_1 + \text{cov}(\beta_1 \beta_2 X_2 + u, X_1 - \beta_2 X_2) / \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

since  $\beta_2$  is constant

$$= \beta_1 + [\beta_1 \beta_2 \text{cov}(X_2, X_1) - \beta_1 \beta_2^2 \text{var}(X_2) - \beta_2 \text{cov}(u, X_2)] / \text{var}(\tilde{X}_1)$$

$$/ \text{var}(\tilde{X}_1)$$

by bilinearity of covariance operator

since  $\text{cov}(u, X_1) = 0$

b.  $\gamma_1 = \beta_1$  iff

$$\beta_1 \beta_2 \text{cov}(X_2, X_1) - \beta_1 \beta_2^2 \text{var}(X_2) - \beta_2 \text{cov}(u, X_2) = 0$$

$$\beta_1 \text{cov}(X_2, X_1) - \beta_1 \beta_2 \text{var}(X_2) - \text{cov}(u, X_2) = 0 \quad \text{or } \beta_2 = 0$$

$$\beta_1 (\text{cov}(X_2, X_1) - \beta_2 \text{var}(X_2)) = \text{cov}(u, X_2) \quad \text{or } \beta_2 = 0$$

$$\beta_1 (\text{cov}(X_2, X_1) - \text{cov}(X_2, \beta_2 X_2)) = \text{cov}(u, X_2)$$

$$\beta_1 (\text{cov}(X_1 - \beta_2 X_2, X_2)) = \text{cov}(u, X_2)$$

$$\beta_1 (\text{cov}(X_1 - \beta_0 - \beta_2 X_2, X_2)) = \text{cov}(u, X_2)$$

$$\beta_1 \text{cov}(X_2, \tilde{X}_1) = \text{cov}(u, X_2)$$

$$\beta_1 \text{cov}(\tilde{X}_1, X_2) = \text{cov}(u, X_2)$$

$$\beta_1 \text{cov}(\tilde{X}_1, X_2) / \text{var}(\tilde{X}_1) = \text{cov}(u, X_2) / \text{var}(X_2)$$

$$\text{cov}(u, X_2) = 0 \text{ (reject) or } \beta_2 = 0$$

$$\beta_2 = \text{cov}(X_1, X_2) / \text{var}(X_2) = 0, X_1 \text{ and } X_2 \text{ uncorrelated}$$



3 Consider the population linear regression model

$$X_3 = \gamma_0 + \gamma_3 W + v$$

$$\text{where } E(v) = 0, \text{ cov}(W, v) = 0$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$$

$$\text{where } E(u) = 0 \text{ and } \text{cov}(X_l, u) = 0 \text{ for } l \in \{1, 2, 3\}$$

$$Y = \beta_0 + \beta_3 \gamma_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 \gamma_3 W + u + \beta_3 v$$

is a population linear regression model of  $Y$  on  $X_1, X_2, W$  such that OLS regression of  $Y$  on  $X_1, X_2, W$  consistently estimates  $\beta_1$  and  $\beta_2$  iff

$$E(u + \beta_3 v) = 0,$$

$$\textcircled{1} E(u + \beta_3 v) = 0$$

$$\textcircled{2} \text{cov}(X_1, u + \beta_3 v) = 0$$

$$\textcircled{3} \text{cov}(X_2, u + \beta_3 v) = 0$$

$$\textcircled{4} \text{cov}(W, u + \beta_3 v) = 0$$

$\textcircled{1}$  holds since

$$E(u + \beta_3 v) = E(u) + \beta_3 E(v) = 0$$

$\textcircled{2}$  holds iff

$$\text{cov}(X_1, u + \beta_3 v) = \text{cov}(X_1, u) + \beta_3 \text{cov}(X_1, v)$$

$$= \text{cov}(X_1, v) = 0$$

$$\text{cov}(X_1, v) = 0$$

By symmetry,

$\textcircled{3}$  holds iff

$$\text{cov}(X_2, v) = 0$$

$\textcircled{4}$  holds iff

$$\text{cov}(W, u + \beta_3 v) = \text{cov}(W, u) + \beta_3 \text{cov}(W, v)$$

$$= \text{cov}(W, v) = 0$$

OLS regression of  $Y$  on  $X_1, X_2, W$  consistently estimates  $\beta_1$  and  $\beta_2$  iff

$$\text{cov}(X_1, v) = \text{cov}(X_2, v) = \text{cov}(W, v) = 0$$

4a All male athletes dropping out of the study does not pose a threat to the internal validity of the study. Assuming that the male athletes decide to drop out of the study for reasons unrelated to the treatment programme, for example, that they imagined it would be more fun to live in the fraternity, the distribution of unmeasured characteristics ~~factor~~ such as ability, and propensity to study, remains identical between the control and treatment group.

Treatment  $X$  remains uncorrelated with unobserved characteristics  $u$ , hence OLS estimators of average treatment effects remain consistent.



b Engineering students sharing a private internet connection poses a threat to the internal validity of the study. This is a case of partial compliance, where engineering students assigned to the control group ~~access the internet in~~ have dorm room internet connections. If this is unknown to the experimenter, OLS estimators underestimate average treatment effects, and are not consistent.

c That art majors in the treatment group never learnt to access their internet accounts does not pose a threat to the internal validity of the study. This is not a case of partial compliance since the art majors in the treatment group still receive the treatment (internet access in their dorm rooms). Unobserved characteristics, including whether an individual would learn to access his internet account if he had <sup>own</sup> internet access in connection in his dorm room, remain uncorrelated with treatment, hence OLS estimators of average treatment effect remain consistent.

d That economics majors in the treatment group provided access to their internet connection to those in the control group for a fee poses a threat to the internal validity of the study. This is a case of partial compliance, where fee-paying members of the control group have dorm room internet connections. If this is unknown to the experimenter, OLS estimators underestimate average treatment effects, and are not ~~consist~~ consistent.

e The damage to the campus network does not pose a threat to the internal validity of the study. Assuming that the rooms affected by the storm are random, or that students are randomly assigned to rooms, and that the affected rooms are known and excluded from the calculation of OLS estimates, the distribution of unobserved characteristics between the treatment group and control group remain identical. OLS estimators of average treatment effect remain consistent.

5. Consider the OLS regression model

$$Y_i = \beta_0 + \beta_1 D_i + \tilde{u}_i$$

where  $E\tilde{u} = 0$  and  $\text{cov}(0, \tilde{u}) = 0$

$$\hat{\beta}_1 = \text{cov}(Y, D) / \text{var}(D)$$

$$= \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})(D_i - \bar{D}) / \frac{1}{n} \sum_{i=1}^n (D_i - \bar{D})^2$$

$$= \sum_{i=1}^n (Y_i - \bar{Y})(D_i - \bar{D}) / \sum_{i=1}^n (D_i - \bar{D})^2$$

$$= \sum_{i=1}^n (Y_i - \bar{Y}) D_i / \sum_{i=1}^n (D_i - \bar{D}) D_i$$

since  $\bar{D}$  is a constant and

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \text{ hence and } \bar{D} = \frac{1}{n} \sum_{i=1}^n D_i$$

$$\sum_{i=1}^n (Y_i - \bar{Y}) = \sum_{i=1}^n (D_i - \bar{D}) = 0$$

$$= \sum_{i=1}^n \sum_{j=1}^n (Y_i - \bar{Y}) \sum_{j=1}^n (D_j - \bar{D})$$

$$= \sum_{i=1}^n \sum_{j=1}^n (Y_i - \bar{Y}) \sum_{j=1}^n (1 - \bar{D})$$

$$= \sum_{i=1}^n \sum_{j=1}^n (Y_i - \bar{Y}) / n_i (1 - \bar{D})$$

since  $1 - \bar{D}$  is constant

$$= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n Y_i - \bar{Y} / (1 - n_i / n_0 + n_1)$$

$$= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n Y_i - \bar{Y} / \frac{n_0}{n_0 + n_1}$$

$$= \frac{n_0 + n_1}{n_0 n_1} \sum_{i=1}^n \sum_{j=1}^n Y_i - \frac{n_0 + n_1}{n_0} \bar{Y}$$

$$= \frac{n_0 + n_1}{n_0 n_1} \sum_{i=1}^n \sum_{j=1}^n Y_i - \frac{n_0 + n_1}{n_0} \frac{1}{n} \sum_{i=1}^n Y_i$$

$$= \frac{n_0 + n_1}{n_0 n_1} \sum_{i=1}^n \sum_{j=1}^n Y_i$$

$$= \frac{n_0 + n_1}{n_0 n_1} \sum_{i=1}^n \sum_{j=1}^n Y_i$$

$$= \frac{1}{n_1} \sum_{i=1}^n \sum_{j=1}^n Y_i - \frac{1}{n_0} \sum_{i=1}^n \sum_{j=1}^n Y_i$$

$$= \frac{1}{n_1} \sum_{i=1}^n \sum_{j=1}^n Y_i - \frac{1}{n_0} \sum_{i=1}^n \sum_{j=1}^n Y_i$$

$$= \frac{1}{n_1} \sum_{i=1}^n \sum_{j=1}^n Y_i - \frac{1}{n_0} \sum_{i=1}^n \sum_{j=1}^n Y_i$$

6. Regression (3) of income transfer on income transfer on observable pretreatment characteristics to test for random receipt of treatment. This entails testing the hypothesis that the coefficients on the controls are zero in regression (3). If treatment is correlated with observable pretreatment characteristics, then assignment of treatment is nonrandom and experimental outcomes reflect both the effect of treatment and the effect of the nonrandom assignment.

Could it be correct to call income transfer the independent variable here? Or is this unclear since income transfer is treated as the dependent variable in the regression?

p-value

$$p = P(F_{6, \infty} > F) = 0.30275$$

$H_0$ : coefficients on all observable treatment characteristics = 0

$H_1$ : coefficient on at least one observable pretreatment characteristic  $\neq 0$

$$F = 1.20, q = 6$$

p-value

$$p = P(F_{6, \infty} > F) = 0.30275$$

Reject the null at all levels of significance  $\alpha \leq 0.30275$ . Fail to reject the null hypothesis that treatment is uncorrelated with



observable pretreatment characteristics at all levels of significance  $\alpha < 0.30275$

- b If the economist regresses food consumption on total household income rather than income transfer, then the coefficient gives the ~~the~~ average change in food consumption associated with ~~a~~ a unit increase in household income. Household income is ~~not~~ ~~not~~ ~~not~~ correlated with other unobserved determinants of food consumption such as nature of breadwinner's occupation. ~~the~~ OLS regression of food consumption on household income does not consistently estimate the effect of a change in ~~total~~ household income on food consumption since orthogonality does not hold. Assuming that treatment is ~~random~~ (successfully) randomly assigned, income transfer is uncorrelated with unobserved determinants of ~~household income~~ <sup>food consumption</sup> and the effect of income transfer on food consumption can be consistently estimated by OLS regression.

- c Yes. Assuming that treatment is (successfully) randomly assigned, OLS regression coefficient on income transfer in a single regression of food consumption on income transfer consistently estimates the causal effect of income transfer on food consumption.

Confidence interval

$$CI = [0.659 - 0.123 c_\alpha, 0.659 + 0.123 c_\alpha]$$

where  ~~$\alpha = 0.025$~~

$$P(-c_\alpha < Z < c_\alpha) = 0.95$$

$$c_\alpha = -\Phi^{-1}(0.025) = 1.96$$

$$= [0.41792, 0.90008]$$

There is a 95% probability that the interval  $[0.41792, 0.90008]$  contains the population regression parameter (coefficient ~~of~~ of income transfer in a population single linear regression of food consumption on income transfer), which gives the share of the additional income spent on food.

- d Including the additional regressors in (2) improves the precision with which the effect of treatment is estimated.

The estimated coefficients in ~~(1) and~~ (1) and (2) are similar, but the standard error of the coefficient in (2) is lower.

Intuitively, the inclusion of other determinants of, or factors correlated with, food consumption



helps to account for some of the variation in food consumption that is not accounted for by ~~the~~ income transfer. Variance of residuals, hence standard error of the coefficient of income transfer decreases.

- e Height is plausibly correlated with food consumption since individuals who consume more food are likely to be taller. Including height in (2) helps to account for some of the variance in food consumption, reducing the variance of residuals, hence reducing the standard error of the coefficient of income transfer, and yielding a more precise estimate of this coefficient.

The coefficient cannot be given a causal interpretation since food consumption and height are likely to be simultaneously determined: taller individuals have higher caloric ~~needs~~ ~~and~~ expenditure and thus must consume more food, individuals who consume more food tend to grow taller as a result.

- f No. The logarithm of the income transfer does not correspond to the ~~proport~~ change in income as a proportion of previous week's income.

Regress ~~ln(income transfer + income in week prior to study)~~ (food consumption / food consumption in prior week) on (income transfer / income in prior week), Age, Education, Household Size, and Height.

- 7a. The estimated coefficient gives ~~the~~ ~~is a~~ ~~consistent estimator~~ gives the causal effect of attending a small ~~class~~ kindergarten class on earnings at age 10.

Given that treatment is (successfully) randomly assigned, other determinants of earnings at age 10 are uncorrelated with whether ~~the~~ the individual (in the study) ~~is~~ attends a small or regular size kindergarten class. As ~~estm~~ regression of Y on D consistently estimates the causal effect of interest.

Is this correct?

Is there any more that can be said?

- b Suppose that Y and X are determined by the following causal models.

$$Y = \beta_0 + \beta_1 D + \beta_2 X + U$$

$$X = \gamma_0 + \gamma_1 D + V$$

Given that ~~D is a~~ treatment is (successfully) randomly assigned, it is ~~not~~ uncorrelated with pretreatment characteristics, including unobserved determinants of  $Y$ , ~~the~~ "collected" in  $u$ . ~~the~~ Hence  $\text{cov}(D, u) = 0$ .

In order that OLS regression consistently estimates  $\beta_0$  and  $\beta_1$ , we require that orthogonality holds, i.e.,  $\text{cov}(D, u) = \text{cov}(X, u) = 0$

$$\begin{aligned}\text{cov}(X, u) &= \text{cov}(\gamma_0 + \gamma_1 D + v, u) \\ &= \text{cov}(\gamma_1 v, u) \\ &\text{by bilinearity of covariance} \\ &\text{since } \gamma_0 \text{ is constant and } \text{cov}(D, u) = 0\end{aligned}$$

It is implausible that  $\text{cov}(v, u) = 0$  since  $\gamma$  income, wealth, and cognitive ability are among the unobserved determinants of  $Y$  and  $X$ , "collected" in  $u$  and  $v$  respectively.

orthogonality does not hold, OLS regression of  $Y$  on  $D$  and  $X$  does not consistently estimate the causal effect of attending a small kindergarten class on future earnings.