

Quantitative Economics Problem Set 3

1a The estimated coefficient on $sibs$ is -0.094 . On average, ~~a man with one more sibling has~~ having one more sibling is associated with 0.094 fewer years of education, holding $meduc$ and $feduc$ constant.

That the estimated coefficient is negative is not entirely surprising. Plausibly, ~~the~~ the families of men with a greater number of siblings shared a finite amount of financial and educational resources between a larger number of children, hence such men received less investment from their parents in their education, and so had fewer total years in education.

On average, having $1 \times 0.094 = 0.094$ more siblings is associated with having 1 fewer year of education.

b On average, a man whose mother had one more year in education, holding constant the number of his siblings and his father's total years in education, had 0.131 more total years of education.

Other causal determinants of $educ$ that are unobserved and collected in the population regression residual include country of residence, area of residence (urban / rural), household income, and household wealth. These determinants are not plausibly uncorrelated with $meduc$. ~~A mother who lived in a developed country and in an urban area with greater access to~~ is more likely to have grown up in such a country and area, hence more likely to have had greater access to education and more likely to have higher total years in education. ~~A mother with high total years in education is also more likely to have higher income, a man belonging to such a household would have had higher household income.~~

Then, the unobserved determinants of $educ$ are ~~not~~ correlated with the observed determinants. Orthogonality fails in the causal model, the population regression parameters do not coincide with the causal model parameters. The sample regression estimates are consistent for the former but not the latter, hence do not provide a reliable estimate for how a mother's educational attainment influences that of her son.

Further, the estimated coefficient provides a reliable estimate of the causal effect of interest only if that causal effect is linear. We have no reason to think so in this case.

$$\begin{aligned} c \quad \hat{E}(educ_A | sibs_A = 0, meduc_A = feduc_A = 12) \\ = 10.36 - 0.094(0) + 0.131(12) + 0.210(12) + \hat{\epsilon}_A \\ = 14.452 \end{aligned}$$

$$\begin{aligned} \hat{E}(educ_B | sibs_B = 0, meduc_B = feduc_B = 16) \\ = 10.36 - 0.094(0) + 0.131(16) + 0.210(16) + \hat{\epsilon}_B \\ = 15.816 \end{aligned}$$

$$\text{Predicted difference} = 15.816 - 14.452 = 1.364$$

An intergenerational ~~regression to the mean~~ mean reversion is observed because the effect of ~~an increase in meduc~~ ~~on~~ higher $meduc$ and higher $feduc$ are, on average, associated with ~~a less than one-to-one~~ ~~pro~~ ~~one-to-one~~ higher $educ$, but in a less than one-to-one proportion.

d Other causal determinants and their correlation with the observed determinants are discussed in (b).

2a Age is a causal determinant of $wage$. Older workers tend to have greater knowledge and experience, hence greater productivity, and higher $wage$. Then, including age in the regression model reduces the magnitude hence the variance of the residuals, hence increases the precision of the estimates of the other regression coefficients.

On average, a worker who is one year older has an hourly $wage$ that is 0.51 dollars higher, holding degree status and gender constant.

The coefficient on age consistently estimates the causal effect of age on hourly $wage$ only if age is uncorrelated with other unobserved determinants of $wage$. This is not plausible. The economic conditions at the time that a person entered the workforce is ~~a~~ one such unobserved determinant of $wage$, and is highly correlated with ~~the~~ age. Then, the ~~the~~ population regression coefficient on ~~the~~ age does not coincide with ~~that of the~~ the causal model coefficient, and the sample regression coefficient is consistent for the former but not the latter.

- b On average, female workers had an hourly wage 3.81 less than that of male workers, holding degree status and age constant.

The negative coefficient constitutes ^{decisive} evidence that being female has a negative causal effect (due to discrimination) on wages only if gender is uncorrelated with other unobserved ~~deter~~ ^{entirely} determinants of wage. This is not plausible. ~~One~~ such determinant is sector of employment and ~~a worker's job sector~~ the nature of work. If it is the case that female workers disproportionately seek employment in sectors or lines of work that pay less, then the negative coefficient on female ~~sector~~ is not a consistent estimator of the ~~effect~~ causal effect of gender on wage, and does not constitute decisive evidence of discrimination.

- c No, degree is not plausibly uncorrelated with the other unobserved determinants of wage. One such determinant is general cognitive ability, another is social capital. Workers with high general cognitive ability are more likely to have a degree because admission into and graduation from a degree programme requires some level of cognitive ability. Workers with high social capital are more likely to have a degree because attending a degree programme is costly and ~~the~~ the ability to afford such a programme is likely correlated with parents' social capital. Then, the population regression (of wage on degree, female, and age) parameters do not coincide with the causal model parameters, and the sample regression coefficients are consistent for the former but not the latter.

Supposing that a worker is not directly employed in some sort of family business, parental income is not plausibly a causal determinant of wage. ~~Then, parental income~~ Parental income is presumably correlated with a young worker's social capital. Then, parental income can be used as a proxy for social capital.

There is no apparent proxy for general cognitive ability apart from test scores, but only supposing that these are not known by employers hence not a causal determinant of wage. It is also necessary to suppose that the above proxies are uncorrelated with degree status. Admittedly, this is not very plausible.

Supposing that such plausible proxies exist, regression of wage on ~~female~~ degree, female, age, and these proxies yields a consistent estimate of the causal effect of degree on wage.

No, the regressors are perfectly multicollinear. An increase (decrease) in study is necessarily associated with a one-to-one decrease (increase) in the sum of sleep, work, and leisure. Then it does not make sense to speak of a causal effect of study because ~~there~~ there is no variation in study as such but only reallocations of time from sleep, work, or leisure to study. It does not make sense to hold sleep, work, and leisure constant while varying study.

- b The given causal model cannot be estimated by OLS regression. The residual in each auxiliary regression is zero because of perfect multicollinearity. For example, in the auxiliary regression $\text{study} = \pi_0 + \pi_1 \text{sleep} + \pi_2 \text{work} + \pi_3 \text{leisure} + \text{study}$, OLS regression yields the parameters $\hat{\pi}_0 = 168$, $\hat{\pi}_1 = \hat{\pi}_2 = \hat{\pi}_3 = -1$, $\text{study}^{\text{sample}} = 0$. Then, the ^{sample} variance of the residual is zero. Hence the OLS regression coefficient which is equal to the sample covariance of score and the residual divided by the variance of the residual is undefined.

- c Omit one of the perfectly multicollinear regressors. Then, the coefficients ~~of~~ on the remaining regressors can be interpreted as the average difference in score associated with a higher ~~value~~ value of that regressor and a lower (in one-to-one proportion) value of the omitted regressor. For example, if study is omitted, the coefficient on sleep is the average difference in score associated with having ~~higher~~ sleep one hour more sleep and one hour less study.

Having omitted one regressor, the ~~each~~ residual in ~~the~~ auxiliary regression is non-zero and has non zero sample variance (supposing there is no perfect multicollinearity between the remaining regressors in the data, which could ~~be the~~ fail to hold if, for example, $\text{work} = 0$ for all observations in the data). Then, the parameters of the ~~population~~ causal model can be estimated (not necessarily consistently) by OLS regression.

a) OLS construction, $\beta_1 = \text{cov}(Y, X) / \text{var}(X)$ and $\beta_1^* = \text{cov}(Y^*, X^*) / \text{var}(X^*)$

$$\begin{aligned}\beta_1^* &= \text{cov}(aY+c, bX) / \text{var}(bX) \\ &= ab \text{cov}(Y, X) / b^2 \text{var}(X) \\ &= \frac{a}{b} \text{cov}(Y, X) / \text{var}(X) \\ &= \frac{a}{b} \beta_1\end{aligned}$$

$$\begin{aligned}\beta_0^* &= \bar{Y}^* - \beta_1^* \bar{X}^* \\ &= \bar{E}(aY+c) - \frac{a}{b} \beta_1 \bar{E}(bX) \\ &= c + a \bar{E}Y - a \beta_1 \bar{E}X \\ &= c + a \beta_0\end{aligned}$$

b) The SER is the square root of the sample variance of the ~~residuals~~ OLS residual.

$$\begin{aligned}\hat{u} &= Y - \hat{\beta}_0 - \hat{\beta}_1 X \\ \hat{u}^* &= Y^* - \hat{\beta}_0^* - \hat{\beta}_1^* X^* \\ &= aY+c - (c+a\beta_0) - \frac{a}{b}\beta_1 bX \\ &= aY - a\beta_0 - a\beta_1 X \\ &= a\hat{u}\end{aligned}$$

$$\begin{aligned}S^2_{\hat{u}} &= \frac{1}{n-k-1} \sum_{i=1}^n (\hat{u}_i - \bar{\hat{u}})^2 \\ &= \frac{1}{n-k-1} \sum_{i=1}^n \hat{E}(\hat{u} - \bar{\hat{u}})^2\end{aligned}$$

$$\begin{aligned}S^2_{\hat{u}^*} &= \frac{1}{n-k-1} \hat{E}(\hat{u}^* - \bar{\hat{u}^*})^2 \\ &= \frac{1}{n-k-1} \hat{E}(a\hat{u} - \bar{a\hat{u}})^2 \\ &= \frac{1}{n-k-1} a^2 \hat{E}(\hat{u} - \bar{\hat{u}})^2 \\ &= a^2 S^2_{\hat{u}}\end{aligned}$$

$$S^*_{\hat{u}^*} = a S_{\hat{u}}$$

The SER is an absolute measure of the variability of Y around fitted values \hat{Y} and so scales in direct proportion to the scaling of Y .

$$\begin{aligned}SSR &= \sum_{i=1}^n \hat{u}_i^2 \\ SSR^* &= \sum_{i=1}^n \hat{u}_i^{*2} \\ &= \sum_{i=1}^n (a\hat{u}_i)^2 \\ &= a^2 SSR\end{aligned}$$

$$\begin{aligned}TSS &= \sum_{i=1}^n (Y_i - \bar{Y})^2 \\ TSS^* &= \sum_{i=1}^n (Y_i^* - \bar{Y}^*)^2 \\ &= \sum_{i=1}^n (aY_i + c - \bar{E}(aY_i + c))^2 \\ &= \sum_{i=1}^n (aY_i + c - a\bar{E}Y_i - c)^2 \\ &= a^2 \sum_{i=1}^n (Y_i - \bar{Y})^2 \\ &= a^2 TSS\end{aligned}$$

$$R^2 = 1 - SSR/SS$$

$$\begin{aligned}R^{*2} &= 1 - SSR^*/TSS^* \\ &= 1 - SSR/SS \\ &= R^2\end{aligned}$$

R^2 is a relative measure of the variability of Y around fitted values \hat{Y} , hence R^2 is invariant to

scaling (and/or translation) of Y and/or X .

8 Let the following denote the OLS regression model of X_1 on X_2

$$\begin{aligned}X_1 &= \pi_0 + \pi_2 X_2 + \tilde{X}_1 \\ \text{where } \bar{E}(\tilde{X}_1) &= \bar{E}(X_2 \tilde{X}_1) = 0\end{aligned}$$

Then, by construction, $\hat{\pi}_0 = \bar{E}X_1 - \hat{\pi}_2 \bar{E}X_2$ and $\hat{\pi}_2 = \text{cov}(X_1, X_2) / \text{var}(X_2)$

By substitution into the OLS regression model of Y on X_1 and X_2

$$\begin{aligned}Y &= \hat{\beta}_0 + \hat{\beta}_1(\pi_0 + \pi_2 X_2 + \tilde{X}_1) + \hat{\beta}_2 X_2 + \hat{u} \\ &= (\hat{\beta}_0 + \hat{\beta}_1 \pi_0) + \hat{\beta}_1 \tilde{X}_1 + (\hat{\beta}_2 + \hat{\beta}_1 \pi_2) X_2 + \hat{u}\end{aligned}$$

where $\bar{E}\tilde{X}_1 \hat{u} = 0$ given that \tilde{X}_1 is a linear function of X_1 and X_2 and each of X_1 and X_2 is orthogonal to \hat{u} , and $\bar{E}X_2 \hat{u} = \bar{E}\tilde{X}_1 \hat{u} = 0$ by construction of the OLS regression model of Y on X_1, X_2 . Then,

$$Y = \hat{\delta}_0 + \hat{\beta}_1 \tilde{X}_1 + \hat{\delta}_2 X_2 + \hat{u}$$

where $\hat{\delta}_0 = \hat{\beta}_0 + \hat{\beta}_1 \pi_0$ and $\hat{\delta}_2 = \hat{\beta}_2 + \hat{\beta}_1 \pi_2$ is a OLS regression model of Y on \tilde{X}_1 and X_2 .

Hence, $\hat{\beta}_1$ satisfies $\hat{\beta}_1 = \text{cov}(Y, \tilde{X}_1) / \text{var}(\tilde{X}_1)$.