

PROBLEM SET 2: ESTIMATING THE RETURNS TO EDUCATION

DUE BY 11:59 PM PDT ON THURSDAY 10/26

A surprisingly large share of the past half-century of research in labor economics has focused on the return to education: the added earnings power that an individual obtains by staying in school an extra year. In the 1970s, the late economist Jacob Mincer formulated what is now seen as the standard relation between human capital and wages:

$$\ln(w_i) = \beta_0 + \beta_1 ed_i + \beta_3 exper_i + \beta_4 exper_i^2 + \varepsilon_i$$

where  $w_i$  is the hourly wage,  $ed_i$  is years of education, and  $exper_i$  is years of labor market experience. This equation is known as the Mincerian Wage Equation. In this problem set, we will explore the difficulties that arise in estimating the returns to education using OLS.

We will use two datasets, both containing data on labor earnings and education among US adults. One is a sample of working-age (25-64) adults in the Current Population Survey, a nationally-representative monthly survey of the non-institutionalized population. This dataset is from March 2018, with data on labor market outcomes in 2017. The other dataset comes from the National Longitudinal Survey of Youth, a study that first surveyed a sample of 14-21 year olds in 1979 and then re-surveyed them annually or biennially to the present. The labor market data are for 2018, when the cohort was aged 53-60.

*To install R and RStudio, follow this link. You are encouraged to work in a group of up to 4 members. You may write code together, but you must write verbal answers yourself. Please use a Markdown template for your code. Write verbal answers in the comments within the Markdown file, so that you produce a single PDF with code, results, and writing, which you will upload to Gradescope.*

1. List your group members.
2. Interpret the Mincerian Wage Equation conceptually in 3-5 sentences. If one assumes that education and experience are exogenous, how should one interpret  $\beta_1$ ? Why do you think the equation has a squared term in experience?
3. Clear the environment, load the relevant packages for your analysis, and load the Current Population Survey ([https://github.com/tvogl/econ121/raw/main/data/cps\\_18.Rdata](https://github.com/tvogl/econ121/raw/main/data/cps_18.Rdata)). Drop anyone who worked fewer than 50 weeks, worked fewer than 35 hours in a typical week, or has 0 dollars in annual earnings. Generate a log hourly wage variable, where the hourly wage equals annual labor earnings divided by

- annual work hours. Generate race/ethnicity dummies for the categories “Black,” “Native,” “Asian,” and “other/multiple.” Generate a new education variable to measure years of schooling. The survey reports education in categories (e.g., less than high school, some high school, etc.), so you will have to decide a reasonable way to convert the categories to years. Generate a “potential experience” variable as follows:  $exper_i = age_i - ed_i - 5$ . Also generate  $exper_i^2$ . Summarize the data. Write 1-3 sentences about any notable patterns in the summary statistics.
4. Estimate the Mincerian Wage Equation using heteroskedasticity-robust standard errors. Interpret in 1-2 sentences. What is the estimated return to education?
  5. Estimate an “extended” Mincerian Wage Equation that controls for race and sex. Does the estimated return to education change after controlling for these covariates? Explain in 1-2 sentences.
  6. In the “extended” regression, is the black-white log wage gap statistically different from the female-male log wage gap? Explain in 1-2 sentences.
  7. Run the “extended” regression separately for men and women. (In `feols()`, you can run separate regressions for each value of `var` by specifying the option `split = ~var`.) By how much do the estimated returns differ by sex? Based on the two sets of regression results, assess whether the difference is statistically significant. Report your findings in 1-2 sentences.
  8. Estimate the male-female difference in returns by adding interaction terms to the “extended” regression in the full sample. Do you get the same answer? (You should.) Instead of the *difference* in returns between men and women, we might instead be interested in the *ratio* of returns. Use the delta method to do inference on the ratio of the returns for women to the return for men. Is the ratio significantly different from 1? Explain in 2-4 sentences.
  9. Now load the NLSY data (<https://github.com/tvogel/econ121/raw/main/data/nlsy79.Rdata>). The NLSY oversampled black and Hispanic/Latino respondents. The variable `perweight` is a sampling weight that can be used to obtain statistics that are representative of the population. Calculate the means of `black` and `hisp` with and without using sampling weights. Which means provide unbiased estimates of the racial/ethnic composition of US adults who were teenagers in 1979? Explain in 2-4 sentences.
  10. Generate a log hourly wage variable and a “potential experience” variable as above. Drop anyone who worked less than full time (e.g., 35 hrs/week for 50 weeks). Estimate an extended Mincerian Wage Equation (controlling for race/ethnicity and sex), with and without using sampling weights. How does the use of sampling weights change the results? Decide whether you want to use sampling weights for the rest of the analysis, and justify your choice. Your answer should be 3-5 sentences overall. For the remainder of the analysis, only use your preferred method.
  11. How do your preferred estimates of the return to education and the return to experience compare to the

estimates from the CPS? If there are differences, hypothesize why. Answer in 3-5 sentences.

12. Does your preferred estimate of  $\beta_1$  represents the causal effect of education? Explain in 2-4 sentences.
13. NLSY respondents took a cognitive test, the Armed Forces Qualifying Test (AFQT), in 1981. They also responded to several questions on their childhood environment. The dataset contains both the cognitive test scores and the measures of the childhood environment. Do you think any of these variables would be appropriate as control variables in the Mincerian Wage Equation? If so, re-estimate the equation, controlling for race/ethnicity, sex, and any other variables as you see appropriate. What happens to the estimated return to education? Interpret any changes you observe. Answer in 3-5 sentences.