

# ECON 695 Problem Set 2

September 10, 2025

**Due: 11:59PM Central Time on Friday, 09/26/2025.** Upload a photocopy of your handwritten solutions, or a PDF generated by LaTeX/lyx to Canvas.

Note 1: I highly recommend using lyx to answer non-coding questions. Lyx is a free document processor that interfaces with LaTeX. You will be able to type and see your equations directly, and export the .lyx file to a PDF.

Note 2: I provide a Jupyter notebook (**PS2\_oaxaca.ipynb**) to guide the coding exercise (question 4). Under “lecture2” on Canvas, I also uploaded a Jupyter notebook for running regressions in Python. You can use it as reference.

## 1 Population Regression vs CEF

In lecture 2, we considered three things:  $y_i$ ,  $E[y_i|x_i]$  and  $x_i'\beta^*$  the "infeasible" or population OLS estimator. We showed that we can write:

$$y_i = x_i'\beta^* + u_i$$

where  $\beta^* = \operatorname{argmin}_b E[(y_i - x_i'b)^2]$

1. Show that  $E[x_i u_i] = 0$ .
2. Show  $\beta^* = E[x_i x_i']^{-1} E[x_i y_i]$ .
3. Show  $E[u_i|x_i] = 0 \Rightarrow E[x_i u_i] = 0$ . Hint: Apply the Law of Iterated Expectations.
4. List two cases where  $E[u_i|x_i] = 0$ . Clearly state your assumption and explain why  $E[u_i|x_i] = 0$  under such assumptions.

## 2 Regressions with Dummies

Suppose we have a sample of size  $N$  made up of two groups, denoted 0 and 1. Let  $D_i$  be a dummy variable with  $D_i = 1$  indicating membership in group 1. Finally, let  $N_0$  and  $N_1$  represent the size of the two subgroups (so  $N = N_0 + N_1$ ).

a) Consider an OLS regression model:

$$y_i = \alpha + \beta D_i + u_i$$

Using the first order conditions for the OLS estimates  $(\hat{\alpha}, \hat{\beta})$ , show that

$$\begin{aligned}\hat{\alpha} &= \frac{1}{N_0} \sum_{i \in 0} y_i, \\ \hat{\alpha} + \hat{\beta} &= \frac{1}{N_1} \sum_{i \in 1} y_i\end{aligned}$$

b) Now consider a case where there are 3 groups: 0, 1, 2 and there are 2 dummies:  $D_{1i} = 1$  if  $i$  is in group 1, and 0 otherwise, and  $D_{2i} = 1$  if  $i$  is in group 2, and 0 otherwise, with the regression model:

$$y_i = \alpha + \beta_1 D_{1i} + \beta_2 D_{2i} + u$$

i) find the first order conditions for minimizing the sum of squared residuals, which define the OLS estimator.

ii) using the FOC show that the OLS estimators for the 3-group model will have

$$\begin{aligned}\hat{\alpha} &= \frac{1}{N_0} \sum_{i \in 0} y_i \\ \hat{\alpha} + \hat{\beta}_1 &= \frac{1}{N_1} \sum_{i \in 1} y_i \\ \hat{\alpha} + \hat{\beta}_2 &= \frac{1}{N_2} \sum_{i \in 2} y_i\end{aligned}$$

### 3 Residualized Regression (Frisch-Waugh)

In lecture 4, we showed that the  $j^{th}$  row of the population regression coefficient  $\beta^*$  from the model

$$y_i = x_i \beta^* + u_i$$

can be obtained by first getting the residual from an auxiliary regression of  $x_{ji}$  on all the other  $x's$ :

$$x_{ji} = x'_{(\sim j)i} \pi + \xi_i$$

then forming:

$$\beta_j^* = E[\xi_i^2]^{-1} E[\xi_i y_i]$$

Suppose that we also “residualized” the dependent variable using an auxiliary regression of  $y_i$  on all the other  $x's$ :

$$y_i = x'_{(\sim j)i} \lambda + \phi_i$$

Show that its also true that:

$$\beta_j^* = E[\xi_i^2]^{-1} E[\xi_i \phi_i]$$

(In other words, we could residualize BOTH  $x_{ji}$  AND  $y_i$  and still get the same right answer).

### 4 Oaxaca Decomposition

The goal of this problem set is to extend the analysis of public-private pay gaps in Lecture 3 to *male workers*.

The file **pay.csv** contains 45,404 observations on male workers, age 30-50, with education of 12 or more years, in the March 2012 and 2013 CPS files. The variables are:

govt=1 if government worker

logwage=log hourly wage based on earnings, weeks worked, and average hours per week last year

educ = years of education

hs = 1 if high school grad (12 years of schooling)

somecoll=1 if education is 13-14 years.

college = 1 if BA (16 years of education)

ma = 1 if masters degree (18 years of schooling)

phd = 1 if phd or LLD or medical degree (20 years of schooling)

age= age in years (range 30-50)

NOTE: note that “somecoll” combines 2 levels of education, 13 and 14 years (13 years are people with some college but no degree, 14 years are people with an AA degree from a community college).

Please create 2 new variables: `educ13=1[educ=13]` and `educ14=1[educ=14]` so you will have a total of 6 possible categories for education: **phd, ma, college, educ14, educ13, hs**.

1. Construct a table like the 1st table (page 7) in Lecture 3. Show the means of all the variables for workers in the private (`govt=0`) and public (`govt=1`) sectors, and the differences in these means.
2. Construct a table like the 2nd table (page 8) in Lecture 3. Show the fractions of private and government workers at each education level, and the mean wages of the two groups in each education level, and the public-private wage gap at each level.
3. Using the notation of Lecture 3, let private sector workers be called “group a” (who are the reference group), let government workers be called “group b”, and define the 6 category variables  $D_{gi}$  for each level of education  $g = 1 \dots 6$ . Let  $x'_i = (D_{1i}, D_{2i}, \dots, D_{6i})$ . Find  $\bar{x}^a$  and  $\bar{x}^b$ .
  - (a) Fit an OLS regression model for wages on  $x_i$  for each group, and estimate the coefficients  $\hat{\beta}^a$  and  $\hat{\beta}^b$ . Verify that these coefficients are the mean log wages of groups  $a$  and  $b$ , respectively, for each level of education.
  - (b) Using your estimates of  $\bar{x}^a$ ,  $\bar{x}^b$ ,  $\hat{\beta}^a$  and  $\hat{\beta}^b$ , construct  $\bar{y}^a = \sum_g \bar{p}_g^a \bar{y}_g^a = (\bar{x}^a)' \hat{\beta}^a$  and  $\bar{y}^b = \sum_g \bar{p}_g^b \bar{y}_g^b = (\bar{x}^b)' \hat{\beta}^b$ . Verify that these match the means of wages for groups  $a$  and  $b$  that you construct directly.
  - (c) Using your estimates of  $\bar{x}^a$ , and  $\hat{\beta}^b$  construct

$$\bar{y}_{counterf}^b = \sum_g \bar{p}_g^a \bar{y}_g^b = (\bar{x}^a)' \hat{\beta}^b$$

Find the adjusted wage gap:  $\bar{y}_{counterf}^b - \bar{y}^a$ . How does this compare to the actual gap  $\bar{y}^b - \bar{y}^a$ ?

- (d) Construct the weight  $w_g = N_g^a / N_g^b$  for education categories  $g = 1 \dots 6$ . Use this to construct  $\bar{y}_{counterf2}^b = \frac{\sum_{i \in b} w_g y_i}{\sum_{i \in b} w_g}$  and verify that:

$$\bar{y}_{counterf}^b = \bar{y}_{counterf2}^b.$$

that is,  $\sum_g \bar{p}_g^a \bar{y}_g^b = \frac{\sum_{i \in b} w_g y_i}{\sum_{i \in b} w_g}$