

Econometrics I, Problem Set 1 (Individual)

Jingwan (Amadea) Luo¹

January 27, 2023

This document answers all questions in the problem set 1 of *Econometrics I* taught by Prof. Maite Cabeza Gutiérrez. Content in blue is for grading. Sidenote is for self-memo and output labels. For empirical analysis, the software in use is MATLAB. Codes are saved for Section Annex, the .m files can be found both accompanied with my submission to TA. Luis Ignacio Menéndez García before 27th Jan., 10:15 a.m. and also my Github Repository², in this repository, all the other related files can be found.

¹ IDEA Graduate Programme in Economics, first year student



²https://github.com/AmadeaLuo/IDEA_Econometrics-I

Question 1

Consider the following figure:

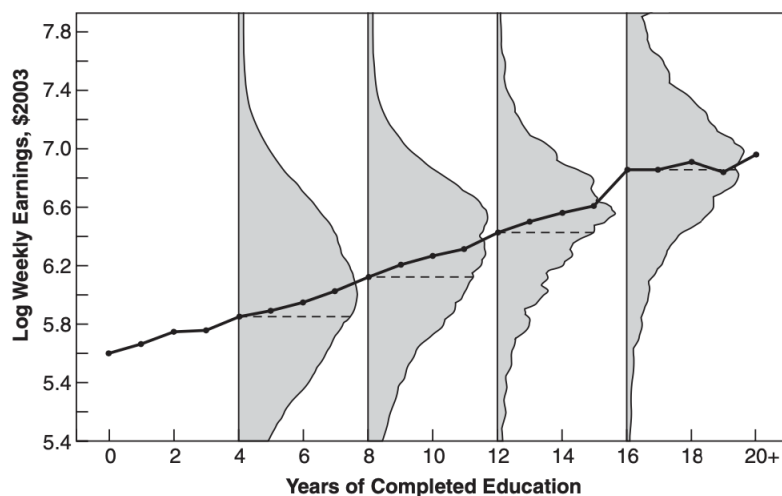


Figure 1: Class Figure 3.1.1 at Slides p.12

Figure 3.1.1 Raw data and the CEF of average log weekly wages given schooling. The sample includes white men aged 40–49 in the 1980 IPUMS 5 percent file.

Figure (1), reproduced from *Mostly Harmless Econometrics: An Empiricist's Companion* (Angrist, 2009), page 31, summarises more than 300,000 observations corresponding to a sample of white males in 1980s in the United States, regarding their years of education (*educ*) and the natural logarithm of their weekly earnings (*lwkywge*).

Datafile `Assig1.csv` includes the observations for variables *educ* and *lwkywge* (i.e. $lwkywge = \ln(wages)$) similar (not identical) to the ones used to produce the figure above.

Question a)

Provide an estimate of average weekly pay for people with 12 years of education. This statistic would be an estimate of which parameter?

Solution.

$$(\overline{wage}|educ = 12) = 397.9337, \quad (a.1)$$

Numerical value computed in Equation (a.1) is an estimate of

$$E(\overline{wage}|educ = 12).$$

*Question b)*

With the help of MATLAB, reproduce the part of Figure 3.1.1 that includes the conditional sample means and the thick black function joining them. (Please notice you are only asked to reproduce the filled circles and the thick black line joining them, not the 4 Kernel density estimates of the conditional distributions included in the plot). Since the data file provided is slightly different than the one used by the authors, you will not get exactly the same values as those included in Figure 3.1.1. Include plot and MATLAB code as your answer.

Solution. Here we give the average log of weekly wage conditional on the years of education. Figure (2) is produced by firstly compute the conditional means for each sub-group of years of education, then joining them with a line. Everything except for the probability distribution on Figure (1) is reproduced.

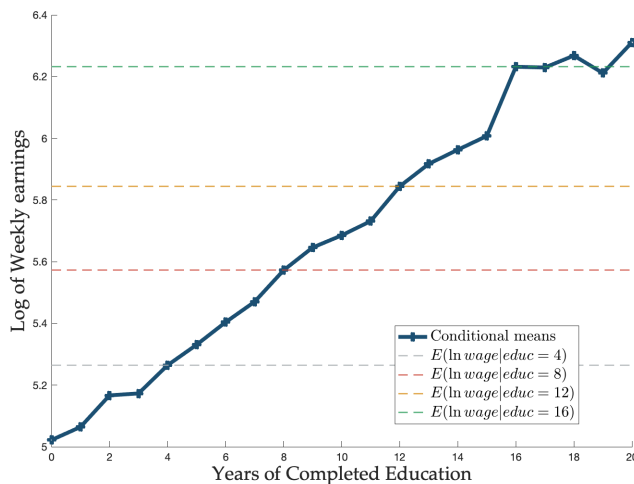


Figure 2: REPLICATION OF FIGURE 3.1.1

Question c)

What does the thick black function you got allow you to say about the relationship between education and earnings? Explain in just a sentence. Try to be as rigorous as possible.

Solution. On *average*, a positive relationship (*not necessarily a causal relationship*) between log of weekly earnings and years of education is observed, since logarithmic transformation is monotone, one can posit that wage has a positive (*not necessarily linear*) relationship with years of the education **in this sample**. ■

Question d)

Consider that now, instead of focusing on conditional sample means, we focussed on conditional sample medians. With the help of MATLAB calculate the conditional medians, and create a new plot including both, the function joining the conditional sample medians together with the function joining conditional sample means. Include the plot in your answer. Comment.

Solution. COMMENT: The conditional medians are higher than the conditional means for every education group. This observation implies that the distribution of log of weekly earnings is slightly left-skewed (negatively skewed).

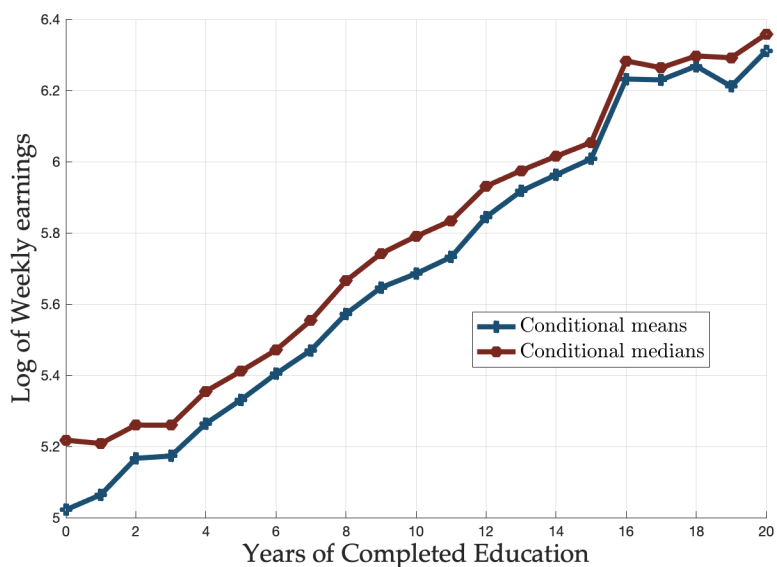


Figure 3: CONDITIONAL MEANS AND MEDIAN FOR LOG OF WEEKLY EARNINGS

We could check this by plotting the density (in Figure (4)).

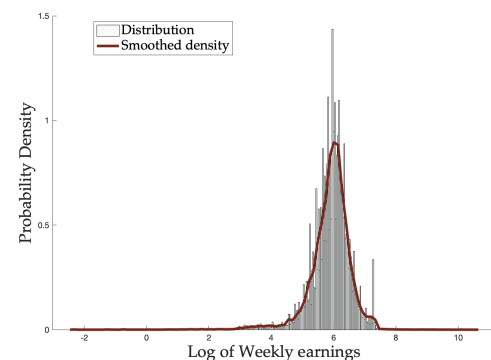


Figure 4: DISTRIBUTION OF LOG OF WEEKLY EARNINGS

Question e)

Let us focus our attention to the thick black function you got in (1b) joining all the conditional sample means (we will ignore now the conditional medians). Your job now is to fit a line through these sample means using MATLAB *polyfit* function.

- + Present a plot including the original sample means and the fitted line.
- + What value did you get for the slope of the line?

Solution. Our goal is to perform a linear fit to the following:

$$\text{lwkwge} = \beta_1 + \beta_2 \text{educ},$$

`p = polyfit(x,y,n)` returns the coefficients for a polynomial $p(x)$ of degree n that is a best fit (in a least-squares sense) for the data in *lwkwge*. The coefficients in p are in descending powers, and the length of p is 2. After the fit, the result is:

$$p = [\beta_2, \beta_1] = [0.0686, 5.0053],$$

Namely,

$$\text{lwkwge} = 5.0053 + 0.0686 \text{educ}.$$

So the slope of the line is $\beta_2 = 0.0686$. Please find the fitted line along with the conditional mean at Figure (5). ■

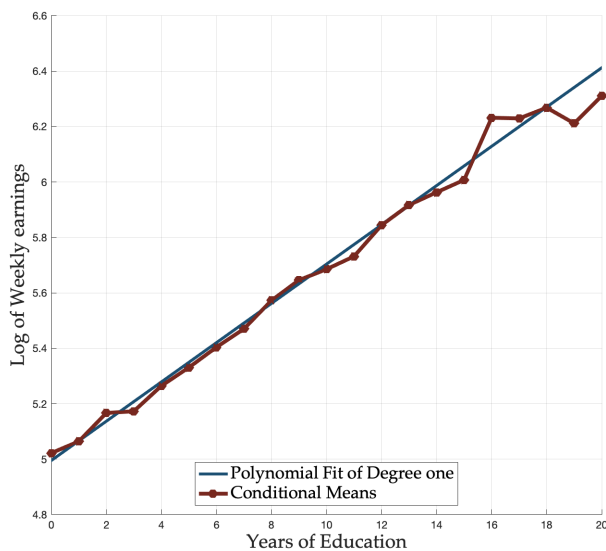


Figure 5: POLYFIT AND CONDITIONAL MEAN

Question f)

Complete this sentence: “According to the value of the slope of the line fitted in the question above, we can say that...” Be as specific as possible.

Solution. ...on **average**, an additional year of education means a 7.09 % raise in weekly earnings for this sample. ■

Question 2

Question a)

Consider the following simple linear regression model:

$$y = \beta_1 + \beta_2 x_2 + \epsilon,$$

where $\epsilon \equiv y - (\beta_1 + \beta_2 x_2)$ and

$$\beta_1 = E(y) - \beta_2 E(x_2) \quad \text{and} \quad \beta_2 = \frac{\text{cov}(x_2, y)}{\text{var}(x_2)}$$

Prove that $\beta_1 + \beta_2 x_2$ is the BLP, in MSE sense, of variable y given variable x_2 .

Solution. In MSE sense,

$$\arg \min_{y_p(x_2)} E\{y - y_p(x_2) \mid x_2\}^2,$$

Expand gives,

$$\arg \min_{y_p(x_2)} E[y^2 \mid x_2] - 2E[y \mid x_2]E[y_p(x_2) \mid x_2] + E[y_p(x_2)^2 \mid x_2]$$

With F.O.C set to 0 then we reached:

$$y_p(x_2) = E(y \mid x_2) \quad (*)$$

which is essentially **what we want to show**.

Below, the **main technique** in use is that “Expectation of a random variable is a constant, and therefore can be taken out of the Expectation expression”.

Consider $E(y \mid x_2)$ and expand y :

$$\begin{aligned} E(\beta_1 + \beta_2 x_2 + \epsilon \mid x_2) &= \beta_1 + \beta_2 x_2 + E(\epsilon \mid x_2) \\ &= \beta_1 + \beta_2 x_2 + E(y - \beta_1 - \beta_2 x_2 \mid x_2) \end{aligned} \quad (2.1)$$

KEY: If substitute $\beta_1 = E(y) - \beta_2 E(x_2)$, $\beta_2 = \frac{\text{cov}(x_2, y)}{\text{var}(x_2)}$ in Equation (2.1) can give us back $E(y \mid x_2)$, then it is indeed the predictor

that Equation (*) requires.

$$\begin{aligned}
 \beta_1 + \beta_2 x_2 + E(y - \beta_1 - \beta_2 x_2 \mid x_2) &= \overbrace{E(y) - \beta_2 E(x_2)}^{\beta_1} + \overbrace{\frac{\text{cov}(x_2, y)}{\text{var}(x_2)}}^{\beta_2} x_2 \\
 &\quad + E(y \mid x_2) - E \left[\underbrace{E(y) - \beta_2 E(x_2)}_{\beta_1} \mid x_2 \right] - \underbrace{E \left(\frac{\text{cov}(x_2, y)}{\text{var}(x_2)} \mid x_2 \right)}_{\equiv A} \cdot x_2; \\
 &= E(y) - \beta_2 E(x_2) + E(y \mid x_2) \\
 &\quad - E \left[\underbrace{E(y)}_{\text{a number}} \mid x_2 \right] + \beta_2 E \left[\underbrace{E(x_2)}_{\text{a number}} \mid x_2 \right]; \\
 &= E(y) - \beta_2 E(x_2) + E(y \mid x_2) - E(y) + \beta_2 E(x_2) \\
 &= E(y \mid x_2).
 \end{aligned}$$

So we reached the requested result.

On a further note, to prove $A = \beta_2 x_2$:

$$\begin{aligned}
 E \left(\frac{\text{cov}(x_2, y)}{\text{var}(x_2)} \mid x_2 \right) &= \frac{E \left\{ \overbrace{E(x_2 y) - E(x_2) E(y)}^{\text{scalars}} \mid x_2 \right\}}{E \left\{ \overbrace{E(x_2)^2 - [E(x)]^2}^{\text{scalars}} \mid x_2 \right\}} \\
 &= \frac{E \{ E(x_2 y) \mid x_2 \} - E \{ E(x_2) E(y) \mid x_2 \}}{E \{ E(x_2^2) \mid x_2 \} - E \{ [E(x)]^2 \mid x_2 \}} \\
 &= \frac{E(x_2 y) - E(x_2) E(y)}{E(x_2^2) - [E(x)]^2} \\
 &= \frac{\text{cov}(x_2, y)}{\text{var}(x_2)}.
 \end{aligned}$$

■

Question b)

Prove the *Decomposition Theorem*, included in slide 1(12), that states that we can always decompose a variable y as:

$$y = E(y|x) + \epsilon$$

with

+ $E(\epsilon|x) = 0$ (ϵ mean independent of x);

+ ϵ uncorrelated with x and with any function of x .

Solution. We know the following equality must hold:

$$y = y + E(y|x) - E(y|x)$$

Write it as

$$y = E(y|x) + \{y - E(y|x)\},$$

and define $\epsilon \equiv y - E(y|x)$.

We check the two properties in turn:

$$[E(\epsilon|x) = 0] : E(y - E(y|x)|x) = E(y|x) - E(y|x) = 0, \quad \checkmark$$

$$[E(x\epsilon) = 0] : E(x\epsilon) = E\left(x \cdot \underbrace{E(\epsilon|x)}_{\text{L.I.E., and } =0}\right) = 0 \quad \checkmark$$

Therefore the decomposition theorem holds. ■

Annex

IMPORTANT Codes can also be found at the submission .zip to TA's email and Github Repository³. Before running the codes, please install the MATLAB add-on "*Professional Plots*" (Aalok, 2023).

³https://github.com/AmadeaLuo/IDEA_Econometrics-I

```

1 %% Setup
2 clc;clear;
3
4 % Set the working directory to the place where the ...
   current file is saved
5 tmp = matlab.desktop.editor.getActive;
6 cd(fileparts(tmp.Filename));
7
8 dt = readtable("Assig1.csv", ...
   'VariableNamingRule','preserve');
```

To Accompany Question 1a)

```

1 %% Question 1 a) estimate of average weekly pay with educ=12
2
3 % +++ Recover variable `wage` from `lwkwlywge`: +++
4 %
5 % Recover by taking exponetial of `lwkwlywge`:
6 wage = exp(dt.lwkwlywge);
7 % Make a table of the recovered wage data:
8 wage1 = array2table(wage);
9 % Glue to a new dataset:
10 dt1 = horzcat(dt, wage1);
11
12 % +++ Conditional mean E(wage | educ = 12): +++
13 %
14 % Create a Sub-sample with educ == 12:
15 sample_12 = dt1(dt1.educ == 12,:);
16 % E(wage | educ = 12):
17 cmean_wage = mean(sample_12.wage);
```


To Accompany Question 1b)

```

1  %% Question 1 b) plot mean of `lwklwyge` conditional on ...
   each educ group
2
3  % +++ Conditional mean: +++
4  %
5  % Get unique educ values
6  educ_groups = unique(dt1.educ);
7  % Initialize vector to store means
8  cmeans = zeros(length(educ_groups),1);
9  % Compute mean for each educ value
10 for i = 1:length(educ_groups)
11     cmeans(i) = mean(dt1.lwklwyge(dt1.educ == ...
        educ_groups(i)));
12 end

```

```

1  % +++ Replicate line chart in Fig 3.1.1: +++
2  %
3  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
4  %           get the add on called           %
5  %           `professional plot`           %
6  %           first                           %
7  %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
8  %
9  % Before anything, set the graph aesthetics
10 PS = PLOT_STANDARDS();

```

```

1  % Plot the means as a line chart with markers
2  figure(1);
3  fig1_comps.fig = gcf;
4  grid off;
5  hold on;
6
7  % Get the conditional mean for educ == 4, 8, 12, 16
8  cmeans4 = cmeans(5,:);
9  cmeans8 = cmeans(9,:);
10 cmeans12 = cmeans(13,:);
11 cmeans16 = cmeans(17,:);

```

```

1 % Drawing
2 % Main graph:
3 fig1_comps.p0 = plot(educ_groups, cmeans);
4 set(fig1_comps.p0, 'Color', PS.Blue5, ...
5     'LineWidth', 4, 'Marker', 'o');
6 % cmeans for educ == 4, 8, 12, 16:
7 fig1_comps.p1 = yline(cmeans4, 'Color', PS.Grey3, ...
8     'LineStyle', '--', LineWidth=1.5);
9 fig1_comps.p2 = yline(cmeans8, 'Color', PS.Red3, ...
10    'LineStyle', '--', LineWidth=1.5);
11 fig1_comps.p3 = yline(cmeans12, 'Color', PS.Orange3, ...
12    'LineStyle', '--', LineWidth=1.5);
13 fig1_comps.p4 = yline(cmeans16, 'Color', PS.Green3, ...
14    'LineStyle', '--', LineWidth=1.5);
15
16 xlabel('Years of Completed Education', ...
17     'FontSize', 20, 'FontName', 'Palatino');
18 ylabel('Log of Weekly earnings', ...
19     'FontSize', 20, 'FontName', 'Palatino');
20 legend('Conditional means', ...
21     '$E(\ln wage \vert educ = 4)$', ...
22     '$E(\ln wage \vert educ = 8)$', ...
23     '$E(\ln wage \vert educ = 12)$', ...
24     '$E(\ln wage \vert educ = 16)$', ...
25     'Location', 'best', ...
26     'Interpreter', 'latex', ...
27     'FontSize', 16)
28
29 hold off;

```

To Accompany Question 1d)

```

1 %% Question 1 d) plot mean&med of `lwkwlywge` conditional ...
  on each educ group
2
3 % +++ Compute the medians +++:
4 %
5 % Initialize vector to store medians
6 cmeds = zeros(length(educ_groups),1);
7 % Compute mean for each educ value
8 for i = 1:length(educ_groups)
9     cmeds(i) = median(dt1.lwkwlywge(dt1.educ == ...
10        educ_groups(i)));
11 end
12 %

```

```

1 % +++ Plot conditional means and medians +++:
2 %
3 % Drawing
4 figure(2);
5 fig1_comps.fig = gcf;
6 grid on;
7 hold on;
8
9 fig2_comps.p0 = plot(educ_groups, cmeans);
10 fig2_comps.p1 = plot(educ_groups, cmeds);
11
12 set(fig2_comps.p0, 'Color', PS.Blue5, ...
13     'LineWidth', 4, 'Marker', 'o');
14 set(fig2_comps.p1, 'Color', PS.Red5, ...
15     'LineWidth', 4, 'Marker', 'Diamond');
16
17 xlabel('Years of Completed Education', ...
18     'FontSize',20, 'FontName','Palatino');
19 ylabel('Log of Weekly earnings', ...
20     'FontSize',20, 'FontName','Palatino');
21 legend('Conditional means', ...
22     'Conditional medians',...
23     'Location','best',...
24     'Interpreter', 'latex', ...
25     'FontSize', 16)
26
27 hold off;
28 %

```

```

1 %% Additional Checkings for Question 1 d)
2
3 % +++ Check distribution of `lwkwlyge`: +++
4 %
5 % Scott's rule for determining optimal bandwidth
6 bw = 3.5*std(dt1.lwkwlyge)*length(dt1.lwkwlyge)^(-1/3);
7 % Drawing
8 figure(3);
9 fig1_comps.fig = gcf;
10 grid off;
11 hold on;
12 histogram(dt1.lwkwlyge, 'Normalization', 'pdf', ...
13     'BinWidth', bw, 'FaceColor',PS.Grey1);
14
15 [density,xi] = ksdensity(dt1.lwkwlyge);
16
17 % Plot density function
18 plot(xi,density,'LineWidth',3.5,'Color',PS.Red5);
19
20 xlabel('Log of Weekly earnings', 'FontSize',22, ...
21     'FontName','Palatino');
22 ylabel('Probability Density', 'FontSize',22, ...
23     'FontName','Palatino');
24 legend('Distribution', 'Smoothed density',...
25     'FontSize',18, 'FontName','Palatino','Location', 'Best');
26
27 hold off;
28 %

```

To Accompany Question 1e)

```

1  %% Question 1 e) polynomial fitting a line between ...
   `lwklwyge` and `educ`
2  % Fit the polynomial of degree one
3  p = polyfit(educ_groups, cmeans, 1);
4
5  % Get the coefficients of the polynomial
6  coefs = p;
7
8  % Generate the fitted values for the independent variable
9  y_fit = polyval(p,dt1.educ);
10
11 % Plot the data and the fitted line
12 figure(4);
13 fig1_comps.fig = gcf;
14 grid on;
15 hold on;
16
17 fig4_comps.p1 = plot(dt1.educ, y_fit);
18 fig4_comps.p2 = plot(educ_groups, cmeans);
19
20 set(fig4_comps.p1, 'Color', PS.Blue5, ...
21     'LineWidth', 3);
22 set(fig4_comps.p2, 'Color', PS.Red5, ...
23     'LineWidth', 3, 'Marker','x');
24
25 xlabel('Years of Education', 'FontSize',20, ...
26     'FontName','Palatino');
27 ylabel('Log of Weekly earnings', 'FontSize',20, ...
28     'FontName','Palatino');
29 legend('Polynomial Fit of Degree one', 'Conditional ...
30     Means',...
31     'FontSize',18, 'FontName','Palatino','Location', 'Best');
32 hold off;

```

References

Aalok, A. (2023), 'Professional plots', https://www.mathworks.com/matlabcentral/fileexchange/100766-professional_plots. [Online; accessed January 22, 2023].

Angrist, Joshua David;Pischke, J.-S. (2009), *Mostly harmless econometrics: an empiricist's companion*, Princeton University Press.
URL: libgen.li/file.php?md5=1e30b583b25afaca5fdda278dccc57f4