

Econometrics: Problem Set #2

Due on April 5, 2024 at 4:00pm

Professor Ben Faber Section 101

Zachary Brandt

Question 1

The Ministry of Commerce in a large country wants to know the causal effect of membership in a local Chamber of Commerce on firm revenues and profits. Firms pay for their membership and they supposedly benefit from the network of information and contacts that the local Chambers of Commerce offer them. But usually, it is only a small minority of all firms that end up paying for the membership.

The Ministry plans to estimate the causal effect of being a member in a local Chamber of Commerce by letting the Ministry's staff estimate the average percentage change in annual firm sales between firms that are members and the rest of the firms that are non-members of their local Chamber of Commerce.

- A) Write down the OLS regression specification that the Ministry's staff could use to implement their analysis described above. Interpret what the intercept and slope coefficients would capture in such a specification.
- B) Using notation from the Potential Outcomes Framework, briefly explain the concept of the Average Treatment Effect (ATE) to the Minister, and how what they plan to estimate in A) relates to this definition.
- C) Referring to the expressions you use in your answer to B), explain why a randomized control trial (RCT) could be useful, and very briefly describe the basics of the RCT design for how the Ministry could set this up.
- D) The Ministry mentions that it has no legal authority to force firms to become members in their local Chambers of Commerce. Using notation from the Potential Outcomes Framework, explain why this information could be important for the interpretation of the results from the RCT relative to the ATE, and how the Ministry should address this concern in the RCT analysis?
- E) The Ministry talked to other economists, and now it is worried about spillover effects on the control group. The staff don't fully understand what the concern is, however. Briefly explain to them the intuition behind this concern, and explain how they could potentially address it when designing the RCT.

Part A

Write down the OLS regression specification that the Ministry's staff could use to implement their analysis described above. Interpret what the intercept and slope coefficients would capture in such a specification.

Solution

The following OLS regression specification is one that the Ministry could implement:

$$\ln(Y_i) = \beta_0 + \beta_1 D_i + u_i$$

where

the subscript i runs over the observations, $i = 1, \dots, n$;

Y_i is the *dependent variable*, annual firm sales

D_i is the *dummy variable*, $D_i = 1$ if the firm is a local Chamber of Commerce member and 0 otherwise

β_0 is the *intercept* of the regression, the population mean value of annual, non-member firm sales

β_1 is the coefficient on D_i , associating a change in D_i by one unit with a $100\beta_1\%$ change in Y_i

$\beta_0 + \beta_1$ is the population mean value of annual, member firm sales

u_i is the *error term*, all the factors responsible for the difference between predicted and observed values

Part B

Using notation from the Potential Outcomes Framework, briefly explain the concept of the Average Treatment Effect (ATE) to the Minister, and how what they plan to estimate in A) relates to this definition.

Solution

One core challenge with evaluating the causal effect of membership in a local Chamber of Commerce on firm revenues and profits is that we cannot observe a firm i in two different states of the world: one where the firm is a member, and one where it is not. Instead, we can compare the mean outcomes of two firm groups (members vs. non-members) to learn about the true, average treatment effect (ATE) of membership.

The ATE is defined as:

$$ATE = E(Y_i(1) - Y_i(0))$$

where $Y_i(1)$ represents the potential outcome (the annual sales) of a firm i if part of a local Chamber of Commerce and $Y_i(0)$ represents the annual sales for the **same** firm if not.

To compare the mean outcomes of the two groups of firms (members vs. non-members of a local Chamber of Commerce), we instead estimate:

$$E(Y_i(1) | X_i = 1) - E(Y_i(0) | X_i = 0)$$

where we condition our observation of Y_i on X_i , with $X_i = 1$ representing if firm i is part of a local Chamber of Commerce, and $X_i = 0$ if it is not. The above is not necessarily equal to the ATE defined earlier because X_i , or membership status, is *not independent* of potential outcomes (annual firm sales). We would have to show first that, in the absence of a local Chamber of Commerce, member and non-member firms would have to be on average identical (there would have to be no selection bias).

Unpacking the above expression we can see the effect of selection bias:

$$\begin{aligned} E(Y_i(1) | D_i = 1) - E(Y_i(0) | D_i = 1) &= \underbrace{E(Y_i(1) | D_i = 1) - E(Y_i(0) | D_i = 1)}_{\text{Average Treatment Effect on the Treated}} \\ &+ \underbrace{E(Y_i(0) | D_i = 1) - E(Y_i(0) | D_i = 0)}_{\text{Selection Bias}} \end{aligned}$$

The first term is the average treatment effect on the treated (ATT), the difference between $Y_i(1)$ and $Y_i(0)$ for the group of firms that are members of a local Chamber of Commerce. The second term estimates the selection bias by comparing $Y_i(0)$, annual sales before being a member, between members, $X_i = 1$, and non-members, $X_i = 0$. $E(Y_i(0) | X_i = 1)$ cannot be observed directly for both of these terms, because we cannot rerun the experiment a second time where the treatment is not applied. The bias term appears because firms that are members of a local Chamber of Commerce may have different outcomes than firms that are not members, *even in the absence of a local Chamber of Commerce*.

Part C

Referring to the expressions you use in your answer to B), explain why a randomized control trial (RCT) could be useful, and very briefly describe the basics of the RCT design for how the Ministry could set this up.

Solution

In a RCT, treatment and control groups are randomly selected, and the average outcomes of these two groups are compared after the treatment. By randomly allocating treatment status, we can use the simple difference in outcomes between treatment and control groups to estimate the ATE. This is because the randomization removes selection bias by ensuring that member firms and non-member firms are comparable in terms of both observable and unobservable characteristics and that any difference between the two groups is not due to systematic differences. This implies that the $ATE = ATT$.

Connecting back to the potential outcomes framework, $D_i = 1$ indicates if firm i is a member of a local Chamber of Commerce (the treatment), and $D_i = 0$ if the firm is not (D_i indicates treatment status). $Y_i(1)$ are the annual sales (the potential outcome) for a firm i if treated (part of a local Chamber or Commerce), and $Y_i(0)$ are the annual sales for the same firm if not.

The ATE is then defined as:

$$ATE = E(Y_i(1) - Y_i(0))$$

If the treatment assignment D_i is assigned randomly and independent of the potential outcomes of the firms ($D_i \perp (Y_i(1), Y_i(0))$), we can compare the mean outcomes of the two groups (members vs. non-members of a local Chamber of Commerce):

$$\begin{aligned} E(Y_i(1) | D_i = 1) - E(Y_i(0) | D_i = 1) &= \underbrace{E(Y_i(1) | D_i = 1) - E(Y_i(0) | D_i = 1)}_{\text{Average Treatment Effect on the Treated}} \\ &+ \underbrace{E(Y_i(0) | D_i = 1) - E(Y_i(0) | D_i = 0)}_{\text{Selection Bias}} \\ &= E(Y_i(1) | D_i = 1) - E(Y_i(0) | D_i = 1) \quad (\text{selection bias is zero}) \\ &= ATT \end{aligned}$$

because the treatment status is independent of outcomes, we can also show:

$$ATT = E(Y_i(1) - Y_i(0)) = ATE$$

This means that if we randomly assign treatment status (membership in a local Chamber of Commerce) and control status (non-membership), we can estimate the ATE by comparing mean sales outcomes across the two groups of observation. We can report the estimated treatment effect of the RCT design from the OLS regression, $\ln(Y_i) = \beta_0 + \beta_1 D_i + u_i$. The β_1 coefficient is the difference in means between the treatment group relative to the control group. Dividing this difference by its standard error produces a t -statistic that allows us to determine whether the observed difference in annual sales between firms part of a local Chamber of Commerce and firms that are not is statistically significant.

Part D

The Ministry mentions that it has no legal authority to force firms to become members in their local Chambers of Commerce. Using notation from the Potential Outcomes Framework, explain why this information could be important for the interpretation of the results from the RCT relative to the ATE, and how the Ministry should address this concern in the RCT analysis?

Solution

In the case where assigned treatment status, D_i , cannot be enforced on firms, T_i represents the effective treatment status of a firm i . Some firms may not comply with their assigned treatment status: firms assigned to become members of a local Chamber of Commerce may opt not to, and firms assigned to be non-members might become members. However, firms' choices to deviate from their assigned treatment are not random, and T_i may be correlated with annual sales (potential outcomes). $T_i \perp (Y_i(0), Y_i(1))$ is no longer true, which means the difference between the mean outcomes of the control and treatment groups may be influenced by selection bias.

However, the initial assignment status, D_i , is still randomly assigned and may have strong predictive power on actual treatment status, T_i . So, to address our original concern, we can use D_i as an instrumental variable for T_i to estimate the causal effect. To be a valid instrumental variable, D_i must satisfy three conditions: instrument relevance, instrument exogeneity, and the instrument exclusion restriction.

Instrument relevance. The instrument needs to be related strongly enough with the endogenous, explanatory variables of interest T_i . So, as long as the assignment protocol is partially followed, then the actual treatment status will be partially determined by the assigned treatment status, making D_i a relevant instrumental variable.

Instrument exogeneity. The instrument must be uncorrelated with any observed variables that also affect Y_i in the error term. That is, $Cor(D_i, u_i) = 0$. Since initial treatment assignment is random, then D_i is distributed independently of u_i , so the instrument is exogenous.

Instrument exclusion restriction. The instrument must not have a direct effect on Y_i except through its relationship with T_i . That is, conditional on T_i , D_i has no effect on Y_i . Again, since the initial treatment assignment is random, D_i is distributed independently of any other variable than T_i that can have an effect on Y_i , satisfying the exclusion restriction.

Therefore, the original random treatment assignment is a valid instrumental variable, and the following is the two-stage least squares regression:

$$\begin{aligned} T_i &= \beta'_0 + \beta'_1 D_i + u'_i \quad (\text{first stage}) \\ \ln(Y_i) &= \beta_0 + \beta_1 \hat{T}_i + u_i \quad (\text{second stage}) \\ \ln(Y_i) &= \beta''_0 + \beta''_1 D_i + u''_i \quad (\text{reduced form}) \end{aligned}$$

In contrast to our original regression, our point estimate from the second stage is: $\hat{\beta}_{IV}$, which gives us a weighted average of the treatment effect β_{1i} , where the weights increase with the degree to which the instrument D_i influences T_i in the first stage.

$$\hat{\beta}_{IV} = \frac{\beta''_1}{\beta'_1} = \frac{E(Y_i|D_i = 1) - E(Y_i|D_i = 0)}{E(T_i|D_i = 1) - E(T_i|D_i = 0)} = \frac{\text{Difference in mean outcomes between groups}}{\text{Difference in fraction treated between groups}}$$

Part E

The Ministry talked to other economists, and now it is worried about spillover effects on the control group. The staff don't fully understand what the concern is, however. Briefly explain to them the intuition behind this concern, and explain how they could potentially address it when designing the RCT.

Solution

The concern behind spillover effects is that, even if we were to isolate the treatment (membership in a local Chamber of Commerce) to firms only in the treatment group and restrict membership for those in the control group, we still have to worry about the effects that the treatment has on the control group. For example, in the case where membership in a local Chamber of Commerce improves firm sales for members, increased sales may contribute to overall economic surplus and increase sales for firms in the control group. This means that the control group is no longer a reference in which nothing has changed.

One solution would be to

Question 2

To answer this question, we will use R and the dataset “Mexico_PS2.csv” that you can download from bCourses. The dataset contains 1153 Mexican municipalities that reported some amount of local tourism activity (measured by local hotel sales) in the year 2000. Write up the answers below in the same document as above. In addition, also attach the complete code that you used to answer the questions. You are not required to export results in LaTeX, it is enough to show your output in R (using either knitR/RMarkdown or just showing the output, e.g. with screenshots).

- A) Visualize a table that lists the number of observations, the mean, the standard deviation, the minimum value and the maximum value for each of the variables in the dataset (if you had difficulties with this in the previous problem set, use the `summary()` command). Briefly describe what we learn from the table about the sample of Mexican municipalities.
- B) Use the data to obtain an OLS point estimate of the effect of the average years of education on the logarithm of local average monthly household incomes. Show your result in a regression table, using the command `stargazer()`. Comment on the interpretation and statistical significance of your results.
- C) List three plausible arguments why the point estimate in b) could be biased upwards or downwards relative to the true causal effect of local tourism activity on monthly household incomes. Be specific in signing the bias for each argument.
- D) Now your GSI suggests that the natural logarithm of kilometer distance between the center of the municipality and the nearest segment of the US-Mexico border could be a valid instrumental variable for the average years of education in a municipality (because incentives for education are higher in regions exporting to the US). List the assumptions that need to hold true for this to be correct.
- E) Verify if the assumption of instrument relevance is satisfied, and export the results into the same regression table that you used above (again using the `stargazer()` command). Comment on the interpretation and statistical significance of your result.
- F) Now estimate the 2nd stage IV (TSLS) point estimate as suggested by your GSI, and show your result in the same regression table you used before (again using the `stargazer()` command). Comment on the interpretation and statistical significance of your result. In reference to your answer to c), is the difference between the OLS and IV point estimates as you expected or rather not?
- G) Now one of your friends suggests that the distance to the US border is likely correlated with other local characteristics that affect local incomes, such as the logarithm of local tourism sales, the logarithm of the average temperature, the logarithm of the average precipitation, and the proportion of indigenous population. Estimate the first-stage regression and the reduced-form regression both before and after you add those variables as additional controls in the two regressions (estimating 4 regression specifications in total). Export your regression results in a new regression table. Comment on the results and what they imply about the validity of the instrumental variable strategy.

Part A

Visualize a table that lists the number of observations, the mean, the standard deviation, the minimum value and the maximum value for each of the variables in the dataset (if you had difficulties with this in the previous problem set, use the `summary()` command). Briefly describe what we learn from the table about the sample of Mexican municipalities.

Solution

Statistic	N	Mean	St. Dev.	Min	Max
year	1,153	2,000.000	0.000	2,000	2,000
municode	1,153	18,168.720	8,045.480	1,001	32,056
inc_m	1,153	2,583.081	1,257.101	299.757	14,792.910
ind_lang	1,153	0.114	0.221	0.000	0.984
educ_years	1,153	7.031	1.354	2.488	11.870
sales_hotel	1,153	27,120.140	201,033.400	1	5,229,616
logtemp	1,153	5.283	0.212	4.658	5.671
logprecip	1,153	4.304	0.562	1.790	5.746
dist_us_km	1,153	696.171	278.301	6.609	1,348.003

From the table, the sample of Mexican municipalities includes 1,153 units of observation. All units were recorded in the year 2000, as the year statistic has no variation across observations in the sample. The units of observation are Mexican municipalities, which, besides the year statistic, also record the municipality code (municode), average monthly income (inc_m), the proportion of the indigenous population (ind_lang), average years of education (educ_years), local hotel sales (sales_hotel), the logarithm of the average temperature (logtemp), the logarithm of the average precipitation (logprecip), and the distance to the US border (dist_us_km).

Average monthly income. Average monthly incomes have a low mean value and a relatively large standard deviation. The distribution could be right skewed as the mean is far closer to the minimum value than the maximum value.

Proportion of the indigenous population. The mean value for the indigenous proportion of Mexican municipalities population stands at 11.4%. However, the standard deviation is very high, suggesting a broad distribution of proportions. Some municipalities have effectively 0% and some are near 100%.

Average years of education. In contrast to the aforementioned statistics, the variation in average years of educational attainment is smaller, suggesting a tighter distribution. Most observations have statistics close to the sample mean, with few outliers.

Local hotel sales. Local hotel sales vary massively. Many municipalities record negligible sales. The distribution appears to be right skewed with a long tail towards higher sales.

Part B

Use the data to obtain an OLS point estimate of the effect of the average years of education on the logarithm of local average monthly household incomes. Show your result in a regression table, using the command `stargazer()`. Comment on the interpretation and statistical significance of your results.

Solution

	<i>Dependent variable:</i>
	log(inc_m)
educ_years	0.183*** (0.007)
Constant	6.485*** (0.047)
Observations	1,153
R ²	0.400
Adjusted R ²	0.399
Residual Std. Error	0.304 (df = 1151)
F Statistic	767.224*** (df = 1; 1151)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Part C

List three plausible arguments why the point estimate in b) could be biased upwards or downwards relative to the true causal effect of local tourism activity on monthly household incomes. Be specific in signing the bias for each argument.

Solution

Part D

Now your GSI suggests that the natural logarithm of kilometer distance between the center of the municipality and the nearest segment of the US-Mexico border could be a valid instrumental variable for the average years of education in a municipality (because incentives for education are higher in regions exporting to the US). List the assumptions that need to hold true for this to be correct.

Solution

Part E

Verify if the assumption of instrument relevance is satisfied, and export the results into the same regression table that you used above (again using the `stargazer()` command). Comment on the interpretation and statistical significance of your result.

Solution

Table 1:

	<i>Dependent variable:</i>
	educ_years
log(dist_us_km)	-0.424*** (0.052)
Constant	9.738*** (0.332)
Observations	1,153
R ²	0.055
Adjusted R ²	0.055
Residual Std. Error	1.317 (df = 1151)
F Statistic	67.597*** (df = 1; 1151)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Part F

Table 2:

	<i>Dependent variable:</i>	
	log(inc_m)	
	<i>OLS</i>	<i>instrumental variable</i>
	(1)	(2)
educ_years	0.183*** (0.007)	0.453*** (0.044)
Constant	6.485*** (0.047)	4.589*** (0.309)
Observations	1,153	1,153
R ²	0.400	−0.465
Adjusted R ²	0.399	−0.466
Residual Std. Error (df = 1151)	0.304	0.475
F Statistic	767.224*** (df = 1; 1151)	

Note:

*p<0.1; **p<0.05; ***p<0.01

Part G

	<i>Dependent variable:</i>			
	educ_years		log(inc_m)	
	(1)	(2)	(3)	(4)
log(dist_us_km)	−0.424*** (0.052)	−0.034 (0.048)	−0.192*** (0.014)	−0.091*** (0.015)
log(sales_hotel)		0.239*** (0.012)		0.042*** (0.004)
logtemp		−0.714*** (0.153)		−0.105** (0.047)
logprecip		−0.432*** (0.072)		−0.124*** (0.022)
ind_lang		−1.246*** (0.150)		−0.439*** (0.046)
Constant	9.738*** (0.332)	11.537*** (0.787)	9.000*** (0.092)	9.233*** (0.241)
Observations	1,153	1,153	1,153	1,153
R ²	0.055	0.435	0.135	0.367
Adjusted R ²	0.055	0.433	0.135	0.365
Residual Std. Error	1.317 (df = 1151)	1.020 (df = 1147)	0.365 (df = 1151)	0.313 (df = 1147)
F Statistic	67.597*** (df = 1)	176.598*** (df = 5)	180.291*** (df = 1)	133.152*** (df = 5)

Note:

*p<0.1; **p<0.05; ***p<0.01

Question 18

Evaluate $\sum_{k=1}^5 k^2$ and $\sum_{k=1}^5 (k-1)^2$.

Question 19

Find the derivative of $f(x) = x^4 + 3x^2 - 2$

Question 6

Evaluate the integrals $\int_0^1 (1-x^2)dx$ and $\int_1^\infty \frac{1}{x^2} dx$.