

Causal Inference

Problem Set 2

Due: April 18th

Problem 1

This problem provides a quick review of the basic potential outcomes notation we discussed in class.

- a) For a binary treatment $D \in \{0, 1\}$, outcome variable Y , and individual units indexed by i , what is the meaning of Y_{1i} and Y_{0i} ? Describe both in words, and choose an example to illustrate.
- b) What is the difference between Y_{1i} and “ Y_i for a unit that actually received the treatment”? Explain the difference using the example you began in part (a).
- c) Define the average treatment effect (ATE) and average treatment effect among the treated (ATT) using potential outcomes notation. Describe in words what each quantity means.
- d) When will the ATT and the ATE be equal to each other? Prove it.
- e) Show formally that $\tau_{ATE(x)} = E[Y_1 - Y_0 | X = x]$ (i.e. a subgroup average treatment effect) is identifiable given random assignment.

Problem 2

For this problem we will use data from Benjamin A. Olken. 2007. “Monitoring Corruption: Evidence from a Field Experiment in Indonesia.” *Journal of Political Economy*. 115: 300-249. The paper and the data set are available on Piazza.

The objective of this experiment was to evaluate two interventions aimed at reducing corruption in road building projects in Indonesian villages. One treatment was audits by engineers; the other was encouraging community participation in monitoring. This problem focuses on the latter intervention, which consisted of inviting villagers to public meetings where project officials accounted for budget expenditures. The main dependent variable is *pct_missing*, a measure of the difference between what the villages claimed they spent on road construction and an independent estimate of what the villages actually spent. Treatment status is indicated by the dummy variable *treat_invite*, which takes a value of 1 if the village received the intervention and 0 if it did not.

The variables in the data set are:

- *pct_missing*: Percent expenditures missing
 - *treat_invite*: Treatment assignment
 - *head_edu*: Village head education
 - *mosques*: Mosques per 1,000
 - *pct_poor*: Percent of households below the poverty line
 - *total_budget*: Total budget (Rp. million) (determined prior to intervention)
- a) Estimate the average treatment effect in this new dataset, using the difference in means estimator.
 - b) Derive an expression for a conservative estimator of the standard error of the above difference-in-means. Why is this estimator conservative?
 - c) Use the data to estimate the standard error you derived in (b).
 - d) Check the covariate balance in this dataset on all covariates (all variables that are not the treatment assignment or the outcome). Decide on sensible balance test statistics and report them in a table. How do the treatment and control group differ?
 - e) Compare the number of treated to untreated units. Comment on the result, and whether it is good or bad in your view.
 - f) Now use regression to estimate the *SATE* (sample average treatment effect). Is this estimate different from the difference-in-means estimate?
 - g) Using your answer from part (c) (not from part (f)), conduct a *t*-test of the null hypothesis that $SATE = 0$. You may use a normal approximation to determine the critical value. Choose an α level and whether you want to conduct a one-sided or two-sided test, based on what you think is sensible.
 - h) Is the conventional estimated standard error of the OLS estimate of the *SATE* (i.e. $\sqrt{\hat{\sigma}_u^2 (X'X)^{-1}_{(2,2)}}$) different than the estimated standard error of the difference-in-means estimate? Why or why not?
 - i) Re-estimate the *SATE* using three additional regression models: (1) one in which you include all pre-treatment covariates as additional linear predictors, (2) another in which you include arbitrary functions of the covariates (polynomials, logs, interactions, etc.) as additional linear predictors, and (3) a third in which you include demeaned versions of the covariates ($X_i - \bar{X}$) as well as the interactions between each of them and the treatment. Report the treatment effect estimates and their robust standard errors. How do these results vary across the regressions?
 - j) Show formally why the variance differs between the controlled estimate of the treatment effect (i.e. the estimate from the regression including pre-treatment covariates) and the uncontrolled

estimate. To do this, use a simplified setup in which you compare the estimator of the variance of $\hat{\beta}_1$ in the model $y_i = \alpha + \beta_1 D_i + \epsilon_{1i}$ to the estimator of the variance of $\hat{\beta}_1$ in the model $y_i = \alpha + \beta_1 D_i + \beta_2 X_{1i} + \epsilon_{2i}$.

Problem 3

- Consider the fictional data set: POdata_2.csv. In these fictional data, we observe an outcome for each unit both under treatment and under control (which, again, is usually impossible in the real world). In addition, we now have a covariate X . This problem uses the data to explore the assignment mechanism and blocking.

Define a treatment vector as the $N \times 1$ vector T that contains each unit's treatment status. Consider the following experimental designs:

1. Each unit i is treated if $Y_{1i} > 100$ and not treated otherwise.
2. Each unit i has probability of receiving treatment, $Pr(T_i = 1) = 0.5$.
3. Treatment is randomly assigned, with $N/2$ units fixed to appear in treatment and control groups.
4. First, each unit is paired with the other unit with an identical X value. Then, within each pair, treatment assignment is randomized: one unit is treated and the other is not treated, with each possible outcome having equal probability.

To answer the following questions, set the seed at 2 in **R**.

- (a) For each of these designs:
 - (i) Compute the number of potential treatment vectors T .
 - (ii) Compute the probability of obtaining a particular treatment vector, and define the probability that any given i receives the treatment.
 - (iii) Implement the treatment assignment mechanism, and estimate $E[Y_i|T_i = 1] - E[Y_i|T_i = 0]$. For each design: Is your estimator unbiased for the true ATE? If not, calculate the bias. How does the realization of the estimate compare to the true ATE using the individual-level potential outcomes?
- (b) Estimate the ATE for design 4 using a bivariate OLS regression and report the conservative estimate for the standard errors. Can you reject the null that the ATE is zero? How can efficiency be improved?

- (c) Suppose that we estimate the ATE with a difference-in-means between treatment and control groups. Design a Monte Carlo study to compare the Mean Square Error (MSE) of this estimator across the four experimental designs (in **R**, set the seed at 2, and draw 1000 samples to conduct the analysis). Rank order the experimental designs by MSE. Offer an explanation for the ranking. In particular, what component of the MSE drives the differences for each design? Why does this component vary across designs? (Hint: Plotting the individual-level potential outcomes, Y_{0i} vs. Y_{1i} , may be informative.)