

# MY457: Problem Set 3 - Difference in Differences

Pedro Torres-Lopez, Michael Ganslmeier, Daniel de Kadt

Fri/15/Mar

This problem set is formative and will not contribute to your final grade. However, it is designed to build and deepen your conceptual understanding of the material and to practice applying the material in R. Using tools and resources such as ChatGPT and Stack Overflow is allowed and encouraged. Bear in mind that you are ultimately responsible for the work you submit, and that these tools often provide incorrect solutions. Make sure that however you use these tools it aligns with your best interests, and enhances your learning in this class.

This problem set must be submitted on Moodle by 5pm on Sat/23/Mar. You must also use the provided `.Rmd` template to produce a `.pdf` with your answers. If your submission is late, is not a `.pdf`, or is not appropriately formatted, you will not receive feedback on your work.

## 1 Concepts

This question reviews some of the concepts covered in class. Mathematical notation can be a useful tool to explain concepts, but it's important that you understand and can explain the concepts clearly and concisely. If you want to support your explanations with mathematical notation, this page provides a tutorial on including mathematical notation in Rmarkdown.

Consider an encouragement design with  $i \in [1, \dots, N]$  units. We have an instrument/encouragement  $Z_i \in [0, 1]$ , but we cannot guarantee compliance with the actual treatment of interest  $D_i \in [0, 1]$ . We are interested in the effect of treatment  $D_i$  on  $Y_i$ .

**1.1** What is the difference between encouragement and treatment? How do these concepts relate the Intent to Treat Effect (ITT) and the Local Average Treatment Effect (LATE)?

**1.2** Given SUTVA, what are the underlying assumptions necessary for the identification of the LATE the above IV setting?

**1.3** Define the compliance types based on the different scenarios that may occur in this setting. Explain in words what each of the types means.

## 2 Simulations

In this question we will use simulated data to test some of our intuitions about instrumental variables. The advantage of using a simulated dataset is that we have explicit control over the data generating process, and know the 'true' answer to any question we pose.

**2.1.** Explain the code below and relate it to an instrumental variables data generating process. Be sure to calculate both the true ATE and the true ITT and include that in your answer.

```
set.seed(123)

n_obs <- 1000
```

```

U <- rbinom(n_obs, 1, .75)

c_type <- ifelse(U == 1,
                sample(1:3, n_obs, prob = c(0.7,0.1,0.2), replace=T),
                sample(1:3, n_obs, prob = c(0.35,0.3,0.35), replace=T))

tau <- ifelse(c_type == 1, 5000, 1000)
tau <- ifelse(c_type == 3, 2500, tau)

Z <- rbinom(n_obs, 1, .5)

D <- ifelse(Z == 1 & c_type == 1, 1, NA)
D <- ifelse(Z == 0 & c_type == 1, 0, D)
D <- ifelse(c_type == 2, 1, D)
D <- ifelse(c_type == 3, 0, D)

Y0 <- rnorm(n_obs, mean = 50000, sd = 2500) + 25000*U
Y1 <- Y0 + tau

Y <- ifelse(D == 1, Y1, Y0)

data <- data.frame(
  cbind(
    Z,
    D,
    Y0,
    Y1,
    Y
  )
)

```

**2.2** Using Ordinary Least Squares (OLS), naively estimate the treatment effect using only  $D$  as a regressor. What do you find? Does this estimator identify the ATE? Why?

**2.3** Repeat the above analysis using only  $Z$  as a regressor. What do you find? What estimand does this estimator identify?

**2.4** Now, estimate the Local Average Treatment Effect (LATE) for the compliers, using the plug-in Wald estimator, considering both  $D$  and  $Z$ . Do you find any differences when you compare this result to your previous estimates? Explain what you find.

**2.5** Using the Two Stage Least Squares (2SLS) estimator, re-estimate the LATE. Do this both manually (using two `lm()` commands) and using the `AER::ivreg` command. How does your result compare to the previous result? How do the two approaches (`lm` and `ivreg`) differ, if at all?

**2.6** (Extra credit): Show that your answers to Questions 2.2, 2.3, and 2.4 were not due to chance. Over at least 1,000 repeat simulations (note: you will want to remove the `set.seed` step), generate a fresh draw of the data and estimate the sample-specific ATE, ITT, and LATE using the three different regression specifications. Calculate the difference between the estimated quantities and what you know to be the true values of these parameters, and store those differences. Finally, produce a histogram that shows the distributions of the differences over your repeated samples, along with their means. What do you conclude?

### 3 Replication

In this section, we will use real-world data to reinforce what we have learned. We will analyse the dataset employed in *Foreign Aid, Human Rights and Democracy Promotion: Evidence from a Natural Experiment*.

Many countries, especially lower- and middle-income countries, are provided foreign aid with the intention that this will improve the living conditions of those in need. This foreign aid might have some spillover effects by encouraging the protection of human rights and entrenching democratic institutions. Does foreign aid improve human rights and democracy?

The authors use instrumental variables to try and studying the effect of foreign aid, measured as overseas development assistance (ODA) on the the CIRI Empowerment Index (CIRI).

**3.1** Explain in your own words the instrumental variable (IV) design that the authors use to answer the research question. Do you have any concerns about the identifying assumptions?

**3.2** Read into R the replication data set (`final.dta`). Each row in the data is a country-year observation, and note that aside from CIRI, ODA, and Colony, all remaining variables are year or country dummy variables. Naively estimate the effect of aid (ODA) on the CIRI Empowerment Index (CIRI), controlling for two-way fixed effects. Be sure to present your results neatly, showing only relevant statistics for ODA and not for all your fixed effects. What do you find?

Hint: to include all variables in a data frame in one regression you can write `y~.` and to later on exclude some variables you can write `y~.-var`.

**3.3** Estimate the first stage, the second stage, and the ITT (also called the reduced form), again controlling for two-way fixed effects. Again, present the results neatly. What do you find?

**3.4** Using the AER package, use the `ivreg` and estimate the LATE. Are your results the same as in 3.3? Why?

**3.5** Estimate the LATE using the plug-in estimator, this time without using any fixed effects. What do you find? Are your results different from before? Why?