# Exercise: Difference-in-Differences

This week, we will simulate a difference-in-differences methodology and see how well we recover the assumed treatment effect.

1. Let's create fake data for Brazil, working at the level of the 5570 municipalities. One characteristic of these municipalities is the presence of oil reserves. Note that this is a *time-invariant* characteristic; if municipality $i$ has oil reserves in year $t$ it also has oil reserves in year $t + 10$. Let's generate a binomial 'oil' variable which gives each of our municipalities a 50% chance of possessing oil.

2. This is a panel dataset so let's make sure that in our dataset we have (i) an indicator for each `municipality` (a number from 1 to 5570), and (ii) an indicator for each time period. We will work with 3 time periods (years, $t = 0$, $t = 1$ and $t = 2$), so create a column for `year` with entries 0, 1 and 2 so that each municipality has three rows in your dataset, one for each year.

3. Now let's simulate potential outcomes - let's say the outcome is 'voter turnout' - for each municipality. We will assume that in general voter turnout is declining over time, and that the presence of oil reserves also reduces turnout. The treatment effect we assume will have a 5% extra voter turnout.

Create variables so that:

$$y_{0,year} = N(60, 5) - 2 * year - 3 * oil$$

$$y_{1,year} = y_{0,year} + 5$$

4. Treatment $D$ is participation by the municipality in a federal government program and occurs between time periods $t = 1$ and $t = 2$. We will assume that only municipalities with oil receive treatment. Make an indicator variable where each municipality with oil is coded as being in the treated group and the rest as control. (*Note* that we are not coding municipalities as treated only in $t = 2$, we are coding for whether they are a treated 'unit' which applies even in $t = 0$, $t = 1$).

5. Now calculate the observed outcome based on the potential outcomes, treatment status *AND time period*.

6. First, let's run the 'naive' **cross-sectional** observational regression of observed outcomes on treatment, using the data from $t = 2$ (when treatment is active and we have both control and treated units). How does the result compare to our simulation assumptions? Why?

7. Next, let's run the 'naive' **before-after** regression of observed outcomes on year, using the data from $t = 1$ and $t = 2$ only for the treated units. How does the result compare to our simulation assumptions? Why?

8. Now, using the data for $t = 1$ and $t = 2$, let's run a basic difference-in-differences regression of the observed outcomes on treatment, year, and the interaction of treatment and year. How do you interpret the results?

9. Our standard errors are wrong here. Cluster your errors by the cross-sectional unit (municipality). (In this case the difference is very small).

10. We can also do a simpler differences-in-differences-in-means estimate without a regression. Create a 2*2 table of average outcomes for the four groups as shown in the table below. Then calculate the differences in the rows and/or the columns, and finally the difference in the differences. Interpret your result.

| Treatment: | D=0 | D=1 |
|---|---|---|
| t=1 | | |

| Treatment: | D=0 | D=1 |
|---|---|---|
| t=2 | | |

11. An assumption of Difference-in-Differences is that there are parallel trends before treatment occurs. Test whether the treated and countrol groups display parallel trends in the outcome variable between time $t = 0$ and $t = 1$. One way to do this is to run exactly the same difference-in-differences regression but excluding time $t = 2$. Interpret your results.

12. Plot a classic difference-in-differences line graph of the average observed outcome, where the x-axis contains the three time periods, the y-axis the average outcome, and there is one line for the treatment group and one for the control group.

13. Finally, let's try to see what estimate we recover when there are **non-parallel trends** produced by **time-varying confounders**. Recreate your dataset but with the following structure of potential outcomes, which only differs in that the falling turnout trend is only present in oil municipalities. (Remember to calculate observed outcomes again). Interpret the results of the difference-in-differences regression this time.

$$y_{0,year} = N(60,5) - 2 * year * oil - 3 * oil$$

$$y_{1,year} = y_{0,year} + 5$$

14. Create the difference-in-differences line graph for this new dataset with time-varying confounders. (The same as in Q12).