# Difference-in-Differences

A practical guide to understand how to measure impact using regression analysis.

JULIANA
VEGA LACORTE

# Contents

Nov. 4, 2020

# The Treatment Group

A treatment group is the group that is anticipated to experience some change after an intervention, program, or project. It can be a group of individuals, households, firms, buildings, etc.

In a policy program, the treatment group is the target population, or "beneficiaries". But since impact evaluation methods can be applied to other areas different from policy programs, we do not necessarily need to talk about a pre-defined target population. We can simply consider an intervention as something that changed, something that is different now and has the potential of affecting the characteristics of our units of study. So the intervention could be a merger, with the treatment group being employees because we think their working conditions would be affected. Or it could be a new medicine, with the treatment group being those patients that will receive the medication. This is actually where the terms intervention and treatment come from: a medical intervention for treating a disease, hence "treatment group" or "treated".

In our case, the "intervention" refers to the green area development or, more specifically, the park construction. It is a change in the urban development that is expected to have some effects on various aspects in the area. By focusing on real estate prices, your treatment units are real estate properties very close to the park. Even though one could think about an area affected by the park project, we need to work with a definition that allows us to evaluate the impact. So the strategy is to choose as the treatment group those properties more likely to be affected by the new park: the ones within a 0-500 meters walking distance.

# The Control Group

Evaluating impact is difficult because we need to uncover causality. We are only interested on the impact (if any) that the park development had on property prices. So we need to eliminate any other plausible explanation for the change in prices observed in the treatment group.

We want to be able to tell the difference between a price change *caused* by the existence of the park versus a price change that *would have ocurred* even if the park had never been built around those properties in the treatment group. But because we only get to observe one scenario for the treatment group (the construction of the park, and price changes), we need to construct a comparison group for the scenario that we do not observe (no park construction, price changes). That comparison group is the control group.

### What is a control group?

A group that shares very similar characteristics with the group that received the intervention (the construction of the park), but was not itself part of the intervention. That is, the control group should be similar to the area in which the park was built, but with no park or new green area development. For example, since close neighboorhoods tend to be similar, it makes sense to choose properties sufficiently close to the new park, but far enough as to not have been impacted by the park development. Under this logic, you are selecting as control group properties that are close to the park area in spatial terms. But to give you another example, it is also possible to consider areas that are "close" in socio-economic terms even if they do not lie geographically next to the treated area.

We can think of the control group as a comparison group that is going to help us deduce what would have happened to the prices of real estate properties if the park had never been built.

Finding a good control group is challenging because we need to create a convincing and reasonable comparison group. In the impact evaluation literature, this is also called finding a good *counterfactual*- finding what would have happened if the park had not existed. In this sense, the price of a property from the treatment group but in absence of the park is its *counterfactual*. This can never be observed.

One of the indications that we have a good control group is that it has some attributes in common with the treatment group at baseline. That is, the groups were similar before the park was built. When judging the similarity between groups we only care about characteristics that may affect price. Our interest is on the impact of the park on real estate prices and, remember, we need to "clean" any observed change in prices caused by the park construction from changes that can be explained by other factors. So when doing comparisons between the control and treatment groups, it is important to select factors/characteristics of real estate properties that are known to influence their prices.

To examine the similarity between control and treatment groups we can look at the average values for selected variables. We would like, ideally, groups that on average had the same characteristics before the park project started. Some researchers report what is called a "balance check": a table that shows how similar or different the groups were before the intervention. Something to keep in mind is that we can only check what we can oberve and measure. But there may still be differences that we do not observe: "unobservables".

For example, I selected some characteristics that could be interesting for your analysis and are shown in Table 1.

If we have variables that are binary (dummy), we use the proportion (or percentage) of cases having the characteristic. Careful here, because for some variables, the mean would not make much sense. For example, if you take the variable *bezirk* that is coded 2 and 7, averaging those values does not have any meaning. But if we transform it to create a dummy coded 1,0 , then the mean value represents the proportion of  1's, and that would be the proportion of properties located in the east.

Table 1. Treatment-Control Descriptive Statistics Before Press Release

| Variable | Treatment 0-500 meter radius | Control 500-1000 meter radius | p-value |
|---|---|---|---|
| Located in East | 68% | 26% | 0.0 |
| New building | 20% | 4% | 0.0 |
| Balcony | 57% | 33% | 0.0 |
| n | 446 | 433 | |

```
// Stata code
-----------------------------------------------------------------
table radius if pre_press==1, c(mean ost)
table radius if pre_press==1, c(mean neub)
table radius if pre_press==1, c(mean balkonbk)
-----------------------------------------------------------------
```

We see from the table that the groups are, on average, very different during the period 2000 until 2005 before the press release. This is only a first indication that is simply telling us to be cautious about any difference in price we observe after the park construction. Any *post-park* difference in price between treatment and control properties could be due to the differences that already existed between the groups before. That is why is not a good idea to use a simple difference as an impact estimator. Instead, you will construct a double-difference, which is the Difference-in-Differences estimator (more on that on the following section).

When presenting these descriptive statistics it may be useful to discuss whether the difference in groups has something to do with the way the data was collected. For example, does your data include all transactions in the 0-500m, or is it a sample taken from a bigger collection of transactions? Or is the difference due to any structural issue in the area? Why are the majority of

transactions in the 0-500 m happening in the east side, while the majority of transactions in the 500-1000m occur in the west side? It would be important to know if this was how the data was collected, or if something happened during the years 2000-2005 that created this distribution.

## *The t-test*

As you can see from Table 1, p-values are also reported. What we are doing is making statistical inferences about the difference between the mean (averages) from two groups. We are in fact using sample means to infer what happens at the population level, because population means are theoretical means we can not observe. So basically we do not want to say that, for example, there is a difference between 0.57 and 0.33 even if we know there is a *mathematical* difference between the values. We want to test if the difference is *statistically* significant. We are thinking in terms of a distribution of values from which we obtain a sample. The idea is to ask how likely it is to observe a difference even if the two groups had the same distribution of values (this is your null hypothesis, that the two groups are the same). If the difference is very large, it becomes more likely that the values come from two different distributions. In other words, the groups are statistically different (with respect to the mean value).

The p-value helps us make a decision about the statistical test. The null hypothesis is the hypothesis of no difference between treatment and control group. Low values of p-value indicate it is "safe" to reject the null hypothesis, thus concluding there is a difference between treatment and control group. In social sciences, a p-value below 0.05 is considered a standard cutoff point to reject the null.

All three p-values in Table 1 are way under 0.05, so we have to say that there is evidence on the statistical difference between the means in treatment and control groups.

| Variable | Treatment 0-500 meter radius | Control 500-1000 meter radius | p-value |
|---|---|---|---|
| Located in East | 68% | 26% | 0.0 |
| New building | 20% | 4% | 0.0 |
| Balcony | 57% | 33% | 0.0 |
| n | 446 | 433 | |

## Assumptions of the t-test

As with most statistical tests, there are some assumptions made about the data that make sure the test works well. For the t-test, we have two basic assumptions: (1) the observations from each sample are independent, and (2) both sample sizes are large, with n > 30. The second is easy to verify, and there is no problem with your data since each group has around 400 units. And I would say you also do not have a problem with the first assumption. The independence assumption is more problematic when you have individuals and you take some measure on the same individuals at two points in time. The two values are clearly related since they refer to the same individual.

There is also an assumption made about the variance of the data. But this assumption does not invalidate the test. It just creates two situations that require different formulas: (1) the two groups being compared have equal variance, or (2) the two groups have different variances. Fortunately, this is an empirical question that can be answered with a test on the equality of variances before conducting the t-test.

# How to calculate impact?

A very innocent way to estimate the effects of the park construction on real estate prices would be to take the difference between properties with near access to the park (treatment gorup) and those without access (control group), after the park opened. This is known as a "single difference", and some researchers might use it because there is no data before the event or intervention they are investigating.

> *Ex-post single diference*:  *difference in outcomes in the treatment group vs control group, with data taken  after the event.*

> Prices *Treatment* group, *After*  -  Prices *Control* group, *After*

Why is this a naïve estimation of the impact?  There might be better real estate properties in our treatment group than in the control group. The difference we observe in prices could be explained by the development of the park, but also because there are key differences between the groups that are driving the difference in prices. And how can we separate those sources of difference? We simply can't.

Without data before intervention there is no way to tell if the control is a good counterfactual. We do not know how different or similar the groups were before.

But with data after and before, it is possible to use the Difference-in-Differences (DiD) approach and get a better estimate of impact.

The DiD combines an after-before approach with a treatment-control group comparison. So it does not only take the difference between treatment and control, but it also considers how the difference changes over time.

> *DiD or double diference*:  *difference in the differences in outcomes before and after the event.*

# Difference-in-Differences

Difference-in-Differences (DiD) is a method commonly used in labor economics to evaluate the effects of labor market policies. Some examples are the effect of training programs on earnings, and the impact of a minimum wage policy on employment. Because it serves as an estimator of causal effects, it is in general a common technique in the impact evaluation literature.
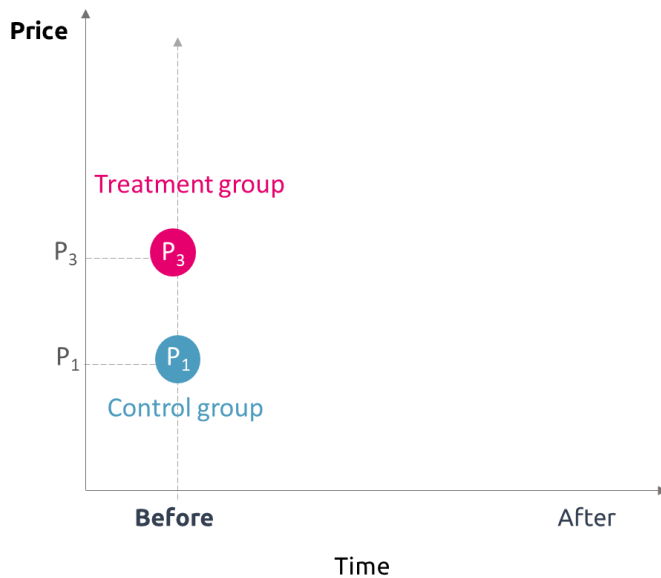
As mentioned briefly before, DiD is a double-difference estimation technique. It is based on the difference in the changes in the outcome between treatment and control groups over time. For your case, the outcome is the average price $\bar{P}$, and the DiD would be:

$$\left( \bar{P}_{Treatment,After} - \bar{P}_{Treatment,Before} \right) - \left( \bar{P}_{Control,After} - \bar{P}_{Control,Before} \right)$$

**Difference** in Ave. Price **over time**, **Treatment** group

**Difference** in Ave. Price **over time**, **Control** group

Taking this double-difference, or DiD approach, will give us a better estimate of the impact that the park development had on property prices. By "better" we mean better than measuring impact using a single difference. The key is observing the difference between treatment and control groups at *both* points in time --before and after the park-- and using it to construct the measure of impact. This should remind you of the counterfactual. In other words, the approach is using the changes in price for properties in the control group to reveal the counterfactual (what would have happened) changes for the treatment properties.

## DiD Graphically

Below you will see a graphical illustration of the DiD estimate. I will take you step-by-step until we get to the final graph showing the DiD estimate.



**Initial Difference**

The average price in the treatment group is $P_3$, while the average price of properties in control group is $P_1$.
(**$P_3$ - $P_1$**) is the difference we observe between the groups *before* the park construction.



**Changes After**

After the park is built, the average price of properties in the control group goes up from $P_1$ to $P_2$ . And the average price in the treatment group goes from $P_3$ to $P_5$.

**Price**

Treatment group

$P_5$

$P_3$
$P_2$

$P_1$

Control group

Before

**After**

Time

**Single-Difference**

After the park is built, the average price of properties in the control group goes up from $P_1$ to $P_2$ . And the average price in the treatment group goes from $P_3$ to $P_5$.
**[$P_5$ - $P_2$]** is the treatment-control difference after the park. That would be a single-difference estimation of impact.

We know that cannot be right. Just comparing average price between treatment and control in the period after is probably not a good estimation of the effect of the park. Average prices were already different between treatment and control properties *before* the park existed. So the difference measured by the segment $P_5 - P_2$ could be due to any underlying differences between the groups (and not because of the park).

Let's turn to the counterfactual. We are able to observe the price increase in properties with near access to the park, after the park opened. But we need to know how much the average price of those properties would have increased without a park, so that we can eliminate it from the impact effect. Here is when the trajectory (change in prices) of the control group comes into play. We will use the change in prices in the control group as an estimation of the change we would have observed in the treatment group without the park.

**DiD**

The difference between the groups would have been $\{P_4 - P_2\}$ in the absence of the park. We substract this difference from the observed treatment-control difference $[P_5 - P_2]$ to obtain the DiD treatment effect $\{P_5 - P_4\}$

With the DiD approach, the true impact of the park is estimated as $(P_5 - P_4)$ because it removes the initial difference between the groups. The ex-post single difference between treatment and control $(P_5 - P_2)$ would have overestimated the effect in the case illustrated here. In empirical cases, it is possible to have either an overestimation or underestimation of the true effect.

What happens if treatment and control group have same values initially? Well, then the single difference and the DiD estimates are equivalent. Think about it as the extreme (and very unlikely) case of having found the perfect counterfactual. If both groups were practically equal ex-ante, then there is no difference to substract from the ex-post observed difference. This case is shown in the next figure.

**Single difference = DiD**

There is no difference between the groups at the beginning (ex-ante). Both have same average price. Then the observed ex-post difference $[P_5 - P_4]$ would be equal to the DiD treatment effect $\{P_5 - P_4\}$

The DiD method works as a good treatment effect estimate thanks to a strong assumption. The assumption that the initial difference between treatment and control would be kept over time until the ex-post period. In other words, we need to assume that both groups experience a common trend in average prices. We are saying that there is no difference in the development of the two real estate markets over time, except for the existence of a park in one of them. This is the parallel trend assumption that I will describe again in a later section.

## DiD Treatment Effect Estimation

One way to compute the DiD is to calculate the sample means of each group-time combination and form the differences. Let's see how it works for the simplest application of DiD estimation: the two-group, two-period case.

This is a general formula, where $\hat{\delta}$ refers to the estimated treatment effect (effect of the policy/intervention), and $\bar{y}$ is the average outcome, like prices por example.

$$\hat{\delta} = \left( \bar{y}_{Treatment,After} - \bar{y}_{Control,After} \right) - \left( \bar{y}_{Treatment,Before} - \bar{y}_{Control,Before} \right)$$

In some texts, you might find the formula a bit different, like this : $\hat{\delta} = \left( \bar{y}_{Treatment,After} - \bar{y}_{Treatment,Before} \right) - \left( \bar{y}_{Control,After} - \bar{y}_{Control,Before} \right)$. But the result is exactly the same.

Let's do some calculations based on your data. Assuming you had data only on two points in time, let's say 1) before the press release in year 2000, and 2) after the opening of the eastpark, year 2012.

Then all we need to do is take the properties in each group and year combination, and calculate their average price. A table like the one below would be created, from which we can obtain the DiD estimation.

### Sample Mean Prices

|  | **Before** Park year=2000 | **After** Park year=2012 | *Difference* |
|---|---|---|---|
| **Treatment** | $\bar{y}_{TB} = 1648.07$ | $\bar{y}_{TA} = 2457.74$ | 809.7 |
| **Control** | $\bar{y}_{CB} = 1669.59$ | $\bar{y}_{CA} = 2873.65$ | 1204.1 |
| *Difference* | -21.52 | -415.91 | -394.39 |

## DiD = -394.39

```
// Stata code
------------------------------------------------------------------
table treat if jahr==2000, c(mean kp_real)
table treat if jahr==2012, c(mean kp_real)
------------------------------------------------------------------
```

This calculation is telling us that average real prices in properties within a 0-500m radius were actually negatively affected by the park. The average price for those properties increased from 2000 to 2012 by 809 eur/m$^2$, but if we believe that properties at the 500-1000m radius make a good control group, we would've expected a larger price increase (by 1204 eur/m$^2$). That is why the DiD in this case is estimating a negative effect of -394 eur/m$^2$.

Note from the table that we can take either the difference for Treatment over time minus the difference for Control over time (809.7 – 1204.1), or the After Treatment vs Control difference - the Before Treatment vs Control difference (-415.91 + 21.52 ) and we will get exactly the same DiD effect estimation.

Now, since you do have more years of data (and not only four points in time) it is also possible to take the average over all those years. Let's see what we would obtain if we divide the before/after period using the opening of the eastpark as a breakpoint. Of course, we are assuming that neither the press release nor the award announcement affected the prices.

| | Before<br>Eastpark Opening<br>2000 - 01.09.2011 | After<br>Eastpark Opening<br>02.09.11-2018 | *Difference* |
|---|---|---|---|
| **Treatment** | 1888.97 | 3299.00 | 1410.0 |
| **Control** | 1930.77 | 3222.26 | 1291.5 |
| *Difference* | -41.8 | 76.7 | 118.54 |

```
// Stata code
--------------------------------------------------------------------------------
table treat if t < date("02sep2011","DMY"), c(mean kp_real)
table treat if t >= date("02sep2011","DMY"), c(mean kp_real)
--------------------------------------------------------------------------------
```

Based on these calculations, we now have an average positive effect on property prices of about 118 eur/m$^2$. The average price of properties near the park increased by 1410 (around 75%). But if we are willing to assume that the average price would have followed the same trend as in the control group in the absence of the park, only around 6.3% increase (118 eur/m$^2$ ) would be attributed to the construction of the park.

There is one drawback in using the sample means this way to calculate the treatment effect. It says nothing about the *statistical significance* of the DiD estimate. For this reason, it is more appropriate to use regression analysis.

## Regression Specification

The DiD approach can also be setup in a regression framework. We can estimate the DiD effect using ordinary least squares regression methods. For the basic case with two groups ( one treatment, one control), and two time periods (one after the event, one before), the regression function looks like this:

$$y_{it} = \alpha + \beta_1 TREAT_i + \beta_2 AFTER_t + \delta(TREAT_i * AFTER_t) + \epsilon_{it} \qquad (1)$$

The coefficient we care about is $\hat{\delta}$, the estimated interaction parameter. And you'll see why this coefficient on the interaction between the indicator for the post-park period (AFTER) and the treatment group (TREAT) will give us the average treatment effect.

A brief reminder about regression analysis:

In a regression model with dependent variable *y*, the expected value of the dependent variable E[*y*] is estimated. The expected value is a statistical concept that refers to the population mean. What we want to do is to get a good estimation of the average value of *y* using our sample data. For example, if we have some data on two variables, *y* and *x*, we can estimate the average value of *y* for different levels of the variable *x*. That is actually a conditional expected value E[ Y | X = x] because it is conditional on x being equal to a specific value.

Going back to the DiD method, in equation (1) the outcome of interest, *y*, would be real estate prices. It has subindex *it* because each price is observed for a specific real estate property *i* and specific date *t*. This regression equation is just looking at certain combinations of means. We are using *dummy* variables (binary variables coded [1, 0]) to pick up the right means. And we only need two dummy variables to form the means of the four groups for the DiD estimation we calculated in the previous section.

$$y_{it} = \alpha + \beta_1 TREAT_i + \beta_2 AFTER_t + \delta(TREAT_i * AFTER_t) + \epsilon_{it}$$

Here the dummy variables are:

TREAT

= 1 if the property is in the treatment group (radius 0-500m),

= 0 if property is in the control group( radius 500-1000m)

AFTER

= 1 if date of transaction is after the park

= 0 if date of transaction is before the park

If we plug in the possible values for TREAT and AFTER, we'll see how we can form the treatment-control combinations for before and after periods. And we we'll see why $\delta$, the interaction term of TREAT and AFTER, is equal to the DiD treatment effect.

$$y_{it} = \alpha + \beta_1 TREAT_i + \beta_2 AFTER_t + \delta(TREAT_i * AFTER_t) + \epsilon_{it}$$

1 treatment 0 control      1 after 0 before

When the values of the indicator variables are:

| TREAT | AFTER | The regression function gives the mean value E[$y_{it}$] for: | through the parameters: |
|---|---|---|---|
| 0 | 0 | Control Before | $\alpha$ |
| 1 | 0 | Treatment Before | $\alpha + \beta_1$ |
| 0 | 1 | Control After | $\alpha + \beta_2$ |
| 1 | 1 | Treatment After | $\alpha + \beta_1 + \beta_2 + \delta$ |

For example:

$$E[y_{it}|TREAT = 1, AFTER = 0] = \alpha + \beta_1 1_i + \beta_2 0_t + \delta(1_i * 0_t) = \alpha + \beta_1$$

Recall our definition for the treatment effect:

$$\hat{\delta} = \left( \bar{y}_{Treatment,After} - \bar{y}_{Control,After} \right) - \left( \bar{y}_{Treatment,Before} - \bar{y}_{Control,Before} \right)$$

Substituting the definition with the regression parameters from the table above, we can see that all but one parameter is eliminated:

$$[ (\alpha + \beta_1 + \beta_2 + \delta ) - (\alpha + \beta_2)] - [ (\alpha + \beta_1) - (\alpha)]$$

$$\alpha + \beta_1 + \beta_2 + \delta - \alpha - \beta_2 - \alpha - \beta_1 + \alpha$$

$$= \delta$$

And that is why if we estimate the regression model (1), we get the **DiD treatment effect** from the estimated coefficient of the **interaction term** Treat*After, $\hat{\delta}$.

Let's confirm empirically that we would indeed get the same result using a simple regression as we did in the previous calculations. We'll replicate the example where the opening of the Eastpark is taken as reference to divide the period in before- and after-park. Again, for simplicity we are assuming that neither the press release nor the award announcement affected the prices.

Calculation using sample averages:

| | **Before** Eastpark Opening 2000 - 01.09.2011 | **After** Eastpark Opening 02.09.11-2018 | *Difference* |
|---|---|---|---|
| **Treatment** | 1888.97 | 3299.00 | 1410.0 |
| **Control** | 1930.77 | 3222.26 | 1291.5 |
| *Difference* | -41.8 | 76.7 | 118.54 |

Calculation using regression:

```
. reg kp_real treat after treat_after
```

| Source | SS | df | MS |
|---|---|---|---|
| Model | 2.3449e+09 | 3 | 781642789 |
| Residual | 6.2893e+09 | 5175 | 1215323.23 |
| Total | 8.6342e+09 | 5178 | 1667482.83 |

| | |
|---|---|
| Number of obs = | 5179 |
| F( 3, 5175) = | 643.16 |
| Prob > F = | 0.0000 |
| R-squared = | 0.2716 |
| Adj R-squared = | 0.2712 |
| Root MSE = | 1102.4 |

| kp_real | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| treat | -41.79561 | 45.69182 | -0.91 | 0.360 | -131.3709 | 47.77965 |
| after | 1291.492 | 43.17545 | 29.91 | 0.000 | 1206.85 | 1376.135 |
| treat_after | 118.5414 | 61.64396 | 1.92 | 0.055 | -2.306793 | 239.3896 |
| _cons | 1930.769 | 31.31917 | 61.65 | 0.000 | 1869.37 | 1992.168 |

Treatment Effect

// Stata code
-------------------------------------------------------------------
gen after = (t >= date("02sep2011","DMY"))
* creating interaction variable
gen treat_after = treat*after

reg kp_real treat after treat_after

*or directly
reg kp_real treat##after
-------------------------------------------------------------------

As you can see, the DiD treatment effect is estimated to be 118.54 eur/m² using the regression model, and the value is the same as what we estimated before with the sample averages. But now you know why this is the case.

*Is the effect significant?*

As mentioned earlier, one of the advantage of using regression analysis is that we can test the statistical significance of the result directly. The output results from the regression already include a t and p-value for the null hypothesis of no effect. Formally, the null hypothesis is: $H_0: \delta = 0$. That is, if the interaction coefficient equals zero then the property prices were not affected by the park construction, on average. The alternative is that the coefficient is different from zero, $H_1: \delta \neq 0$. We are not specifying a particular direction of the effect, we are only testing if there is *any* effect,

either positive or negative. We would like to find evidence in our data to reject the null of no effect and conclude that there is an effect.

To see if the effect is significant, we look at the p-value of the corresponding coefficient (the one from the interaction: treat_after). Remember, we reject the null thypothesis if the p-value is low enough. The standard threshold in social sciences, is to reject the null at the 5% significance level, if the p-value < 0.05.

| kp_real | Coef. | Std. Err. | t | P>|t| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| treat | -41.79561 | 45.69182 | -0.91 | 0.360 | -131.3709 | 47.77965 |
| after | 1291.492 | 43.17545 | 29.91 | 0.000 | 1206.85 | 1376.135 |
| treat_after | 118.5414 | 61.64396 | 1.92 | 0.055 | -2.306793 | 239.3896 |
| _cons | 1930.769 | 31.31917 | 61.65 | 0.000 | 1869.37 | 1992.168 |

From the results, we see that the p-value = 0.055, and we might be very tempting to change the significance level a little more and reject the null. In theory, the researcher should pick the significance level before doing the test, and not change the selected rule afterwards. However, there are indeed many papers that have also opted for simply reporting the levels and say something like: "the null hypothesis is rejected at the 5.5% significance level", or "the effect was found to be statistically significant at the 5.5% level". The idea behind is that they leave the judgement to the readers. Another common practice is to report three significance levels: 10%, 5%, and 1%. The results table then has the coefficients with asteriks with a note that goes like "* denotes significance at the 10% level; ** denotes significance at the 5% level; *** denotes significance at the 1% level.

## Control Variables

We know that in the DiD framework, the comparison group is known as the control group. Its purpose is to help us deduce what would have happened to the prices of real estate properties if the park had never been built. It does that because it is very similar to the treatment group (we created it that way) except for some characteristics, some of which we may not observe. The DiD uses the difference between the groups before and after the park. So as long as the difference

between treatment and control is constant over time, we can say we are removing differences that we cannot observe. Because we knew they were different before (measured by the difference in average price), we are just assuming the reason for that difference stayed constant until after the park.

The DiD eliminates the effects of factors that do not change over time (and that can affect price) from the impact estimate. In the literature, they refer to this as a model that controls for "unobserved, time-invariant" differences (or heterogeneity). Just keep in mind, this is only true if the unobserved difference between groups is assumed to be constant over time.

Now, a regression model can also "control for" differences in *observed* characteristics that *change* over time. But of course we need 1) data on those characteristics and 2) data on those characteristics for different periods of time. You do have data like this in your dataset. You have data on factors that can affect the price per sqm, and you have multiple time periods.

When we add into a regression model other factors that help account for an explanation to the change in prices that is not related to the intervention (the park construction), we called them "control variables". They are just there to "control" for any other changes in the buildings or apartments in your sample that are unrelated to the park construction.

With control variables, the regression equation looks like this:

$$y_{it} = \alpha + \beta_1 TREAT_i + \beta_2 AFTER_t + \delta(TREAT_i * AFTER_t) + \gamma X_i + \epsilon_{it}$$

where X is just representing all the variables different from *treat* and *after*. The notation is used specially if there are many variables included, otherwise you can show them explicitly in the equation. For example, let's say you decide to include location (whether the property is in the west or east) and an indicator for new building as control variables. Your equation would then look like this:

$$y_{it} = \alpha + \beta_1 TREAT_i + \beta_2 AFTER_t + \delta(TREAT_i * AFTER_t) + \gamma_1 ost_i + \gamma_2 neub_i + \epsilon_{it}$$

Running the regression in Stata gives you the output shown in the next page.

## DiD regression with control variables

| Source | SS | df | MS |
|--------|-----|------|------------|
| Model | 5.0166e+09 | 5 | 1.0033e+09 |
| Residual | 3.6176e+09 | 5173 | 699328.137 |
| Total | 8.6342e+09 | 5178 | 1667482.83 |

| Number of obs = | 5179 |
|---|---|
| F( 5, 5173) = | 1434.69 |
| Prob > F = | 0.0000 |
| R-squared = | 0.5810 |
| Adj R-squared = | 0.5806 |
| Root MSE = | 836.26 |

| kp_real | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---------|----------|-----------|-------|-------|-----------|-----------|
| treat | -182.108 | 35.82025 | -5.08 | 0.000 | -252.3309 | -111.8852 |
| after | 1125.549 | 32.86139 | 34.25 | 0.000 | 1061.126 | 1189.971 |
| treat_after | 45.70574 | 46.78563 | 0.98 | 0.329 | -46.01388 | 137.4254 |
| ost | 318.7183 | 25.46744 | 12.51 | 0.000 | 268.7914 | 368.6453 |
| neub | 1641.903 | 28.4966 | 57.62 | 0.000 | 1586.037 | 1697.768 |
| _cons | 1557.825 | 25.75326 | 60.49 | 0.000 | 1507.338 | 1608.312 |

control variables

// Stata code
----------------------------------------------------------------------
reg kp_real treat after treat_after ost neub
----------------------------------------------------------------------

We are using the same specification as before where we defined *AFTER* as the period after the opening of the Eastpark. If we compare these results with control variables to the results without control variables, we'll notice that the estimated treatment effect has changed. Now the construction of the park has no statistically significant effect on the average price. The p-value of the coefficient for the interaction term is 0.329 and thus we cannot reject the null of no effect.

It seems like, once we control for the effect that east-kreuzberg and new building has on prices, the effect of the park goes away. By looking at the coefficients, *neub* suggests there is some large correlation between prices and new building. The same for *ost*. On average, prices tend to be higher for properties located in East-Kreuzberg, and for new buildings. And our data either cannot distinguish those effects from the park effect, or the park really did not affected average prices and the increase observed in the treatment group is due to other factors like new building.

When we are doing a DiD analysis, we usually do not care about interpreting the coefficients on control variables. Even if the coefficients are not sifnificant, they belong to the model if it helps

isolate the effect of interest. We are simply interested on estimating the *treatment effect*, so we only care about the coefficient for the interaction term.

Note on the dependent variable: When the dependent variable is defined to be the natural log, for example *y = log(price),* is because we are interested in the *percentage* increase (or change) in prices. Since this is how you should specify your regressions, in further examples I will use the log(price) as the dependent variable.

## Regression with Multiple Events

So far, we have been working with a model specification that covers the case of having one event that divides the sample in two periods: one before and one after. And that is accomplished with one dummy variable *AFTER*.

To expand the model for more than one event we simply create a specific *AFTER* dummy for each event, and add them to the model interacted with the variable that indicates the group *TREAT*.

A regression equation with three before/after splits looks like this:

$$y_{it} = \alpha + \beta_1 TREAT_i + \beta_2 AFTER1_i + \beta_3 AFTER2_i + \beta_4 AFTER3_i + \delta_1(TREAT_i * AFTER1_i)$$
$$+ \delta_2(TREAT_i * AFTER2_i) + \delta_3(TREAT_i * AFTER3_i) + +\epsilon_{it}$$

Using your dataset, if we create four dummies to split the samples as *post_press*, *post_award*, *post_open1*, and *post_open2*, we get the results shown in the next page.

## DiD regression with four events and control variables

reg lnkp treat##post_press treat##post_award treat##post_open1 treat##post_open2 ost neub

| lnkp | Coef. | Std. Err. | t | P>|t| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| 1.treat | -.0717071 | .0249757 | -2.87 | 0.004 | -.12067 | -.0227442 |
| 1.post_press | .106943 | .0408188 | 2.62 | 0.009 | .026921 | .1869651 |
| treat#post_press | | | | | | |
| 1 1 | -.0342093 | .0618735 | -0.55 | 0.580 | -.1555075 | .0870888 |
| 1.post_award | .0908946 | .0393726 | 2.31 | 0.021 | .0137078 | .1680815 |
| treat#post_award | | | | | | |
| 1 1 | .0790195 | .0603144 | 1.31 | 0.190 | -.0392223 | .1972613 |
| 1.post_open1 | .2863808 | .0218596 | 13.10 | 0.000 | .2435268 | .3292348 |
| treat#post_open1 | | | | | | |
| 1 1 | -.1054873 | .0320622 | -3.29 | 0.001 | -.1683427 | -.0426319 |
| 1.post_open2 | .2227187 | .0208769 | 10.67 | 0.000 | .1817911 | .2636463 |
| treat#post_open2 | | | | | | |
| 1 1 | .0645259 | .0297532 | 2.17 | 0.030 | .0061971 | .1228547 |
| ost | .1624264 | .0111524 | 14.56 | 0.000 | .140563 | .1842898 |
| neub | .5840128 | .012461 | 46.87 | 0.000 | .5595839 | .6084417 |
| _cons | 7.154634 | .0176904 | 404.44 | 0.000 | 7.119953 | 7.189315 |

We care about the interaction terms. According to the results, the coefficients for *treat\*post_press*, and *treat\*post_award* are not significant, while *treat\*post_open1* and *treat\*post_open2* are significant. This suggests that the effect of the park on prices did not appear until it was opened.

## Regression with Time effects

We have extended the DiD regression model with variables that control for characteristics that can affect property prices such as new building, and east/west location. A further extension consists on controlling for changes in average prices that occur simply because time goes by. Sometimes we would like to control for a financial crisis that may have affected prices in a particular year. To do this, we add *time dummies*, which are indicator variables that take the value of 1 if the observation belongs to a specific year and 0 if it is any other year.

For example, a dummy *yr00* will have the value of 1 for data collected in the year 2000, and 0 in the rest of the years. You can create one dummy for each year in your dataset.

| | Year dummies | | | |
|---|---|---|---|---|
| | yr00 | yr01 | ... | yr10 |
| **2000** | 1 | 0 | | 0 |
| **2001** | 0 | 1 | | 0 |
| 2002 | 0 | 0 | | 0 |
| 2003 | 0 | 0 | | 0 |
| 2004 | 0 | 0 | | 0 |
| 2005 | 0 | 0 | | 0 |
| 2006 | 0 | 0 | | 0 |
| 2007 | 0 | 0 | | 0 |
| 2008 | 0 | 0 | | 0 |
| 2009 | 0 | 0 | | 0 |
| **2010** | 0 | 0 | | 1 |
| 2011 | 0 | 0 | | 0 |
| 2012 | 0 | 0 | | 0 |
| 2013 | 0 | 0 | | 0 |
| 2014 | 0 | 0 | | 0 |
| 2015 | 0 | 0 | | 0 |
| 2016 | 0 | 0 | | 0 |
| 2017 | 0 | 0 | | 0 |
| 2018 | 0 | 0 | ... | 0 |

Time dummies in the regression will control for changes in average prices over time that are common to both treatment and control properties. You will have 19 time dummies in total. The thing to watch here is not to include those 19 dummies in the regression. One should always include one dummy less than the total number of categories, so 19-1=18 here. The year dummy that is not included serves as reference category to interpret the rest. Each estimated coefficient tells you if averages prices in that year are higher or lower than your reference year. In Stata, if you use the **i.** notation directly on the regression, you don't have to worry about dropping one dummy out. Stata does that automatically, and it usually takes the first time period as reference.

## DiD regression with time effects (four events and control variables)

| lnkp | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| 1.treat | -.05804 | .0246311 | -2.36 | 0.018 | -.1063274 | -.0097527 |
| 1.post_press | .190723 | .0545494 | 3.50 | 0.000 | .083783 | .2976629 |
| treat#post_press | | | | | | |
| 1 1 | -.0408528 | .0604275 | -0.68 | 0.499 | -.1593163 | .0776108 |
| 1.post_award | -.1945431 | .0601919 | -3.23 | 0.001 | -.3125447 | -.0765414 |
| treat#post_award | | | | | | |
| 1 1 | .0672809 | .0587939 | 1.14 | 0.253 | -.0479801 | .1825419 |
| 1.post_open1 | .1284323 | .0348091 | 3.69 | 0.000 | .0601918 | .1966729 |
| treat#post_open1 | | | | | | |
| 1 1 | -.1136317 | .0315239 | -3.60 | 0.000 | -.1754318 | -.0518315 |
| 1.post_open2 | -.027082 | .0372764 | -0.73 | 0.468 | -.1001595 | .0459956 |
| treat#post_open2 | | | | | | |
| 1 1 | .1040558 | .0294189 | 3.54 | 0.000 | .0463822 | .1617294 |
| ost | .1544225 | .0109262 | 14.13 | 0.000 | .1330025 | .1758425 |
| neub | .6040269 | .0125792 | 48.02 | 0.000 | .5793663 | .6286874 |
| jahr | | | | | | |
| 2001 | -.1190235 | .0439116 | -2.71 | 0.007 | -.2051088 | -.0329381 |
| 2002 | -.1313805 | .0429693 | -3.06 | 0.002 | -.2156187 | -.0471423 |
| 2003 | -.2294251 | .0442179 | -5.19 | 0.000 | -.316111 | -.1427392 |
| 2004 | -.189774 | .0392197 | -4.84 | 0.000 | -.2666613 | -.1128867 |
| 2005 | -.2219786 | .0445356 | -4.98 | 0.000 | -.3092873 | -.13467 |
| 2006 | -.2412847 | .0705774 | -3.42 | 0.001 | -.3796463 | -.1029231 |
| 2007 | -.0261099 | .0849845 | -0.31 | 0.759 | -.1927156 | .1404959 |
| 2008 | .0441883 | .0836422 | 0.53 | 0.597 | -.119786 | .2081625 |
| 2009 | .0330912 | .0834233 | 0.40 | 0.692 | -.1304539 | .1966363 |
| 2010 | .1167965 | .0819659 | 1.42 | 0.154 | -.0438915 | .2774845 |
| 2011 | .1515214 | .082801 | 1.83 | 0.067 | -.0108038 | .3138466 |
| 2012 | .2198501 | .0875706 | 2.51 | 0.012 | .0481746 | .3915256 |
| 2013 | .3003783 | .0902033 | 3.33 | 0.001 | .1235415 | .4772151 |
| 2014 | .3831663 | .093703 | 4.09 | 0.000 | .1994685 | .5668641 |
| 2015 | .4108809 | .093884 | 4.38 | 0.000 | .2268283 | .5949335 |
| 2016 | .4937458 | .0938842 | 5.26 | 0.000 | .3096929 | .6777987 |
| 2017 | .5555171 | .0940889 | 5.90 | 0.000 | .3710628 | .7399713 |
| 2018 | .6691479 | .0958287 | 6.98 | 0.000 | .4812828 | .8570129 |
| _cons | 7.301102 | .0322377 | 226.48 | 0.000 | 7.237902 | 7.364301 |

// Stata code

--------------------------------------------------------------------------------------------------------------

reg lnkp treat##post_press treat##post_award treat##post_open1 treat##post_open2 i.jahr

--------------------------------------------------------------------------------------------------------------
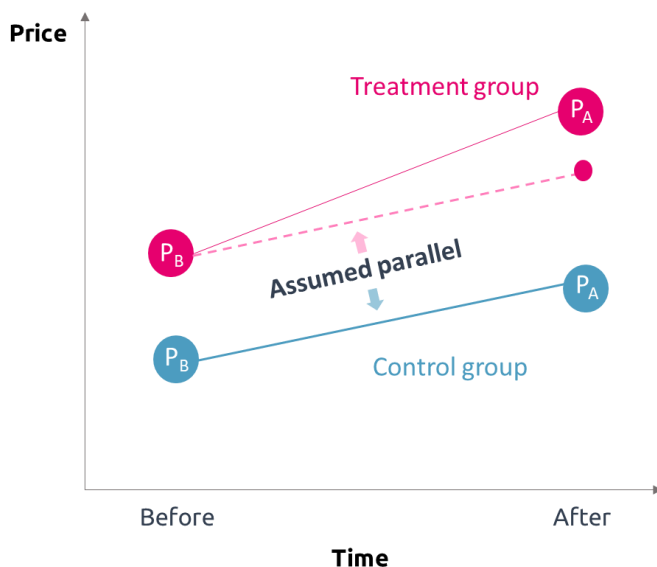
Just as with control variables, we do not need to interpret each time dummy coefficient. But as an example, the coefficient on the year 2001 means that prices in 2001 were, on average, around 12% lower than prices in year 2000.

The interaction terms that measure the impact of the park are giving us similar results as before. Both the press release and award announcement are not significant, while the opening of the park are statistically significant. The effect of opening the west park is estimated to have increased prices by 10% on average.

# Parallel Trend Assumption

DiD estimation relies on a strong assumption known as the parallel trends assumption. The parallel trend assumption implies that real estate prices follow the same trend over time in both treatment and control groups when no park development takes place.

In the figure below you can notice how the trend line for the control group is used to create the counterfactual trend for the treatment group. That is, we assume that the price development in the treatment group would have been parallel to the price development of properties in our control group.



The trend line observed in the control group is assumed to be the trend the treatment group would have followed if the park had not been built.

The dashed line represents how we imagine the prices would have grown in the abscence of the park. It is parallel to the blue line, so both lines have the same slope. That is, both lines (blue and dashed pink) represent the same rate of change.
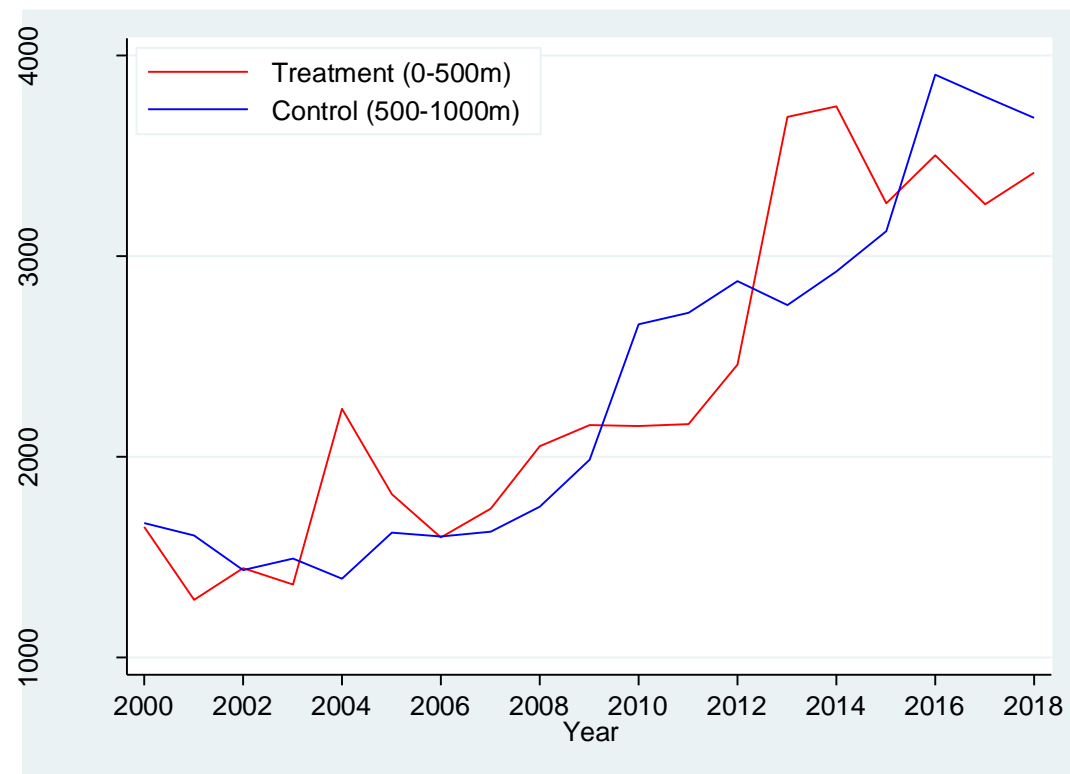
The idea is that if the treatment and control groups were trending similarly before the construction of the park, then they likely would have continued to trend similarly from, let's say 2011-2018, if the park had not been built. The DiD relies on this assumption to produce unbiased estimates of the treatment effect.

We can test the assumption using pre-park data. Looking at several periods before, we would like to show that the treatment and control groups have a similar pattern in pre-park years. This gives support to the assumption, but is not a test to demonstrate that the assumption is valid. We simply cannot test the validity of the assumption because we can never observe what would have happened to the treatment in the absence of the park. So any test is only an indirect way of

evaluating how likely it is that the assumption holds. If you provide convincing evidence that the assumption holds, your effect estimates become more credible.

We can use graphs to check the parallel trends assumption in an exploratory step of our analysis. If you plot average prices over time, divided by treatment and control group, you would obtain a graph like this one:

Average Price over Time, by Treatment and Control group.



```
// Stata code
--------------------------------------------------------------------------------------------------------------------------
* generate Mean price per Year and Treat group
bys jahr treat: egen mean_kpreal = mean(kp_real)
* Line plot over time
twoway (line mean_kpreal jahr if treat==1, lcolor(red))   ///
          (line mean_kpreal jahr if treat==0, lcolor(blue)), ///
          legend(label(1 Treatment (0-500m)) label(2 Control (500-1000m)) ///
          region(lstyle(none)) cols(1) ring(0) position(11)) ///
          ytitle("Mean Price per sqm") xtitle("Year") ///
          xlabel(#10)
--------------------------------------------------------------------------------------------------------------------------
```

By visual inspection, we see both groups showing the same downward trend on average prices in the first years of your sample. From around 2003 until 2018 there seems to be an upward tred. But in the context of the parallel trends assumption we care about the pre-treatment period, that is, the period before the existence of the park. For example, let's say we would like to consider the opening of the eastpark as breakpoint for the analysis, then we would like to observe the same trend in the groups before 2011.

In Stata, we can use `lfit` to estimate a linear trend and add it to our graph. The line trend is created by fitting values from a liner prediction of the mean price on time. Just as an approximation let's draw a trend line for the years 2000-2010, considering it as the pre-park period.

In the figure below we can see that both groups follow an upward trend. However, the lines do not look parallel, except for the first years. The slope of the fitted line for the control group seems steeper, so that average prices increased more rapidly than those in the treatment group during that period. What is obvious is that, the ave. price of properties in the treatment group had a big jump in 2004. It would be interesting to know if something in particular caused that change in price.

Trend in years 2000-2010

```
// Stata code
-------------------------------------------------------------------------------------------------------------------------------------------------------------------
twoway (line mean_kpreal jahr if treat==1, lcolor(red))   ///
        (line mean_kpreal jahr if treat==0, lcolor(blue)) ///
        (lfit mean_kpreal jahr if treat==1 & jahr<2011, lcolor(black) clwidth(medthick) ) ///
        (lfit mean_kpreal jahr if treat==0 & jahr<2011, lcolor(black) clwidth(medthick) clpattern(dash)), ///
                legend(label(1 Treatment (0-500m)) label(2 Control (500-1000m)) label(3 trend treatment)
label(4 trend control) ///
                region(lstyle(none)) cols(1) ring(0) position(11)) ///
                ytitle("Mean Price per sqm") xtitle("Year") ///
                xlabel(#10)
```

## Regression with Year*Treat Interactions

There is a regression model that can be used to evaluate the parallel trend assumption. We can interact (multiply) the indicator for the treatment group *TREAT* with a full set of year dummies. That is, we will have interaction terms between each year dummy and the treatment dummy. The model would look like this:

$$y_{it} = \alpha_0 + \alpha_1 TREAT_i + \varphi_1 yr2001_i + \varphi_2 yr2002_i + \cdots + \varphi_{18} yr2018_i + \beta_1(TREAT_i * yr2001_i)$$
$$+ \beta_2(TREAT_i * yr2002_i) + .. + \beta_{18}(TREAT_i * yr2018_i) + X_i\gamma + \epsilon_{it}$$

Here I'm trying to use a similar notation as your supervisor so you can read his model more easily. The $\varphi$ coefficients refer to the time dummies. We are leaving the year 2000 out, so this becomes the base/reference year. Remember, time dummies control for changes in average prices over time that are common to both treatment and control groups.

The $\beta$ coefficients are the ones for the interaction terms. And your control variables are represented by the vector X with corresponding coefficients $\gamma$.

The $\beta$ coefficients on the interaction terms give the estimated difference between the treatment and control group relative to the baseline year. In a way, it is like estimating a treatment effect for each year. Ideally, one would like to observe insignificant coefficients in the pre-treatment period

(years 2001 up to maybe 2011) cause this would mean there is no difference in price between the groups in that period so they must be "trending similar". We can check this by looking at the p-values of the coefficients in the years previous to the park. Remember that they refer to the test that the coefficient is equal to zero (no effect).

Since we have interactions for each year, we can also look at the coefficients on the years after the park, and see if there is an effect there (significant coefficients). We could also look at the magnitude of the coefficients and see whether the effect goes away, stays constant, or increases over time.

In the next page you can see the results for the regression including the two control variables ost and neub. In this regression I am already using robust standard errors that are clustered by id.

## Regression with treat*year interactions for each year

```
Linear regression                                Number of obs =      5179
                                                 F( 39,  5178) =    265.65
                                                 Prob > F       =    0.0000
                                                 R-squared      =    0.5891
                                                 Root MSE       =    .35234

                            (Std. Err. adjusted for     5179 clusters in id)
```
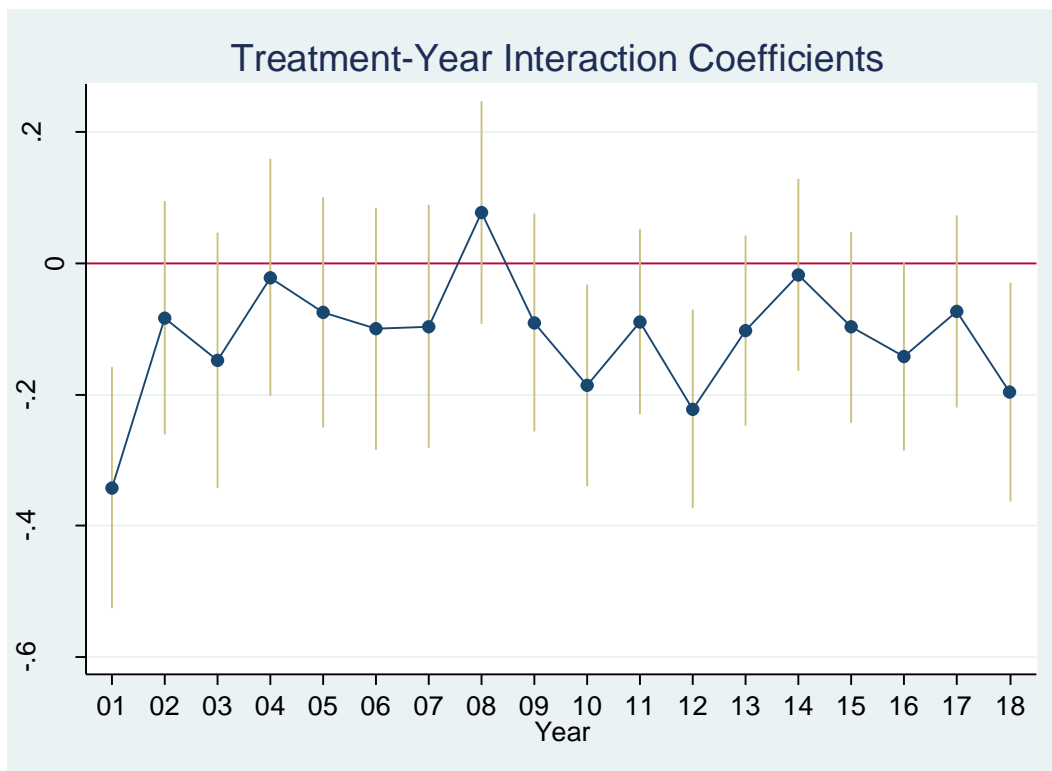
| lnkp | Coef. | Robust Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| 1.treat | .0471178 | .0668377 | 0.70 | 0.481 | -.0839123 | .1781478 |
| **jahr** | | | | | | |
| 2001 | .045249 | .0624478 | 0.72 | 0.469 | -.077175 | .1676731 |
| 2002 | -.1043396 | .0625318 | -1.67 | 0.095 | -.2269284 | .0182492 |
| 2003 | -.1628716 | .0716424 | -2.27 | 0.023 | -.303321 | -.0224223 |
| 2004 | -.193948 | .0689358 | -2.81 | 0.005 | -.3290913 | -.0588046 |
| 2005 | -.1206286 | .0594183 | -2.03 | 0.042 | -.2371136 | -.0041437 |
| 2006 | -.120344 | .0620834 | -1.94 | 0.053 | -.2420537 | .0013657 |
| 2007 | .0166534 | .0599685 | 0.28 | 0.781 | -.1009101 | .134217 |
| 2008 | -.0003292 | .057365 | -0.01 | 0.995 | -.1127889 | .1121305 |
| 2009 | .0752556 | .0562914 | 1.34 | 0.181 | -.0350993 | .1856104 |
| 2010 | .1965329 | .0501347 | 3.92 | 0.000 | .0982477 | .294818 |
| 2011 | .2301089 | .0466666 | 4.93 | 0.000 | .1386226 | .3215951 |
| 2012 | .3979125 | .0481979 | 8.26 | 0.000 | .3034243 | .4924007 |
| 2013 | .4348331 | .0484576 | 8.97 | 0.000 | .3398357 | .5298305 |
| 2014 | .4758528 | .0507431 | 9.38 | 0.000 | .3763748 | .5753307 |
| 2015 | .5537677 | .0485723 | 11.40 | 0.000 | .4585454 | .64899 |
| 2016 | .6528917 | .0452161 | 14.44 | 0.000 | .5642491 | .7415343 |
| 2017 | .6851482 | .0486725 | 14.08 | 0.000 | .5897296 | .7805669 |
| 2018 | .8551224 | .0522439 | 16.37 | 0.000 | .7527023 | .9575426 |
| **treat#jahr** | | | | | | |
| 1 2001 | -.3419204 | .0938402 | -3.64 | 0.000 | -.5258868 | -.157954 |
| 1 2002 | -.0834705 | .0907119 | -0.92 | 0.358 | -.2613042 | .0943631 |
| 1 2003 | -.1478385 | .0991114 | -1.49 | 0.136 | -.3421387 | .0464616 |
| 1 2004 | -.0214502 | .091987 | -0.23 | 0.816 | -.2017836 | .1588832 |
| 1 2005 | -.0748888 | .0893832 | -0.84 | 0.402 | -.2501176 | .1003401 |
| 1 2006 | -.0990841 | .0939993 | -1.05 | 0.292 | -.2833624 | .0851943 |
| 1 2007 | -.0963926 | .0944612 | -1.02 | 0.308 | -.2815766 | .0887913 |
| 1 2008 | .0770336 | .086748 | 0.89 | 0.375 | -.0930291 | .2470964 |
| 1 2009 | -.0901321 | .0845416 | -1.07 | 0.286 | -.2558693 | .0756052 |
| 1 2010 | -.1856409 | .0783472 | -2.37 | 0.018 | -.3392345 | -.0320474 |
| 1 2011 | -.088757 | .0722041 | -1.23 | 0.219 | -.2303075 | .0527934 |
| 1 2012 | -.2224734 | .0772512 | -2.88 | 0.004 | -.3739183 | -.0710285 |
| 1 2013 | -.1027632 | .0737963 | -1.39 | 0.164 | -.2474352 | .0419088 |
| 1 2014 | -.0177685 | .0748519 | -0.24 | 0.812 | -.1645099 | .1289728 |
| 1 2015 | -.0974515 | .0744291 | -1.31 | 0.190 | -.243364 | .048461 |
| 1 2016 | -.1422398 | .0733357 | -1.94 | 0.052 | -.2860087 | .0015291 |
| 1 2017 | -.0728544 | .074581 | -0.98 | 0.329 | -.2190647 | .073356 |
| 1 2018 | -.1963288 | .0849201 | -2.31 | 0.021 | -.3628081 | -.0298495 |
| ost | .1519704 | .0110015 | 13.81 | 0.000 | .1304029 | .173538 |
| neub | .5977827 | .0100785 | 59.31 | 0.000 | .5780247 | .6175408 |
| _cons | 7.263307 | .0423388 | 171.55 | 0.000 | 7.180305 | 7.346308 |

reg lnkp treat##i.jahr ost neub, vce(cluster id)
--------------------------------------------------------------------

## Plotting Interaction Coefficients

In order to see more clearly how the estimated coefficients are changing over time, we can plot them in a graph like the one below.

The graph shows the point estimates along with confidence intervals that we obtained from the regression. We would like to see a flat trend in the years previous to the park. That is, we would like to see a pre–treatment pattern that moves around zero or stays constant. For the regression that I am running here, this does not seem to be the case. The effects change quite a lot in the beginning, and only in the years 2005-2007 we see that they stay the same.

```
// Stata code
---------------------------------------------------------------------------------------------------------------------------------------------------
reg lnkp treat##i.jahr ost neub, vce(cluster id)

global coefinter coefinter 1.treat#2001.jahr 1.treat#2002.jahr 1.treat#2003.jahr 1.treat#2004.jahr ///
        1.treat#2005.jahr 1.treat#2006.jahr 1.treat#2007.jahr 1.treat#2008.jahr 1.treat#2009.jahr ///
        1.treat#2010.jahr 1.treat#2011.jahr 1.treat#2012.jahr 1.treat#2013.jahr 1.treat#2014.jahr ///
        1.treat#2015.jahr 1.treat#2016.jahr 1.treat#2017.jahr 1.treat#2018.jahr

coefplot, recast(connected) vertical keep($coefinter) yline(0) ///
        coeflabels(1.treat#2001.jahr = "01" 1.treat#2002.jahr = "02" 1.treat#2003.jahr = "03" ///
        1.treat#2004.jahr = "04" 1.treat#2005.jahr = "05" 1.treat#2006.jahr = "06" ///
        1.treat#2007.jahr = "07" 1.treat#2008.jahr = "08" 1.treat#2009.jahr = "09" ///
        1.treat#2010.jahr = "10" 1.treat#2011.jahr = "11" 1.treat#2012.jahr = "12" ///
        1.treat#2013.jahr = "13" 1.treat#2014.jahr = "14" 1.treat#2015.jahr = "15" ///
        1.treat#2016.jahr = "16" 1.treat#2017.jahr = "17" 1.treat#2018.jahr = "18") ///
        xtitle("Year") title("Treatment-Year Interaction Coefficients") ///
        ciopts(recast(rcap) pstyle(p8))
```