

Linear Regression Review (How to Interpret Interactions)

Professor Song Yao
Olin Business School

Customer Analytics

Interpreting Interaction Term: One Continuous and One Dummy Variables

Let's consider the determinants of executive pay as an example for interpreting interaction effects

Data: A sample of 1000 executives

Questions:

1. How does work experience (years) affect one's salary?
2. Does the salary differ for men and women?
3. Does the effect of work experience on salary differ for men and women?
-> Is there an **interaction** between experience and gender?

```
# load libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
import statsmodels.formula.api as smf
from tabulate import tabulate

# load salary data
data_url = 'https://tinyurl.com/salary-experience-csv'
salary = pd.read_csv(data_url)
print(salary.head())
```

	experience	female	mba	salary	gender
0	1	0	0	100272.3	Male
1	2	1	1	173588.3	Female
2	3	1	1	127221.0	Female
3	4	1	1	186438.0	Female
4	5	0	1	188325.3	Male

3

Question 1: How does work experience affect one's salary?

$$\text{Salary} = a + b * \text{Experience}$$

```
# linear regression of salary on experience
salary_experience = smf.ols('salary ~ experience', data=salary)
salary_experience_fit = salary_experience.fit()
print(salary_experience_fit.summary())
```

```
=====
                        ...
=====
```

	coef	std err	t	P> t	[0.025	0.975]
<i>a</i> -> Intercept	1.277e+05	3150.212	40.537	0.000	1.22e+05	1.34e+05
<i>b</i> -> experience	5192.5235	178.231	29.134	0.000	4842.772	5542.275

```
=====
```



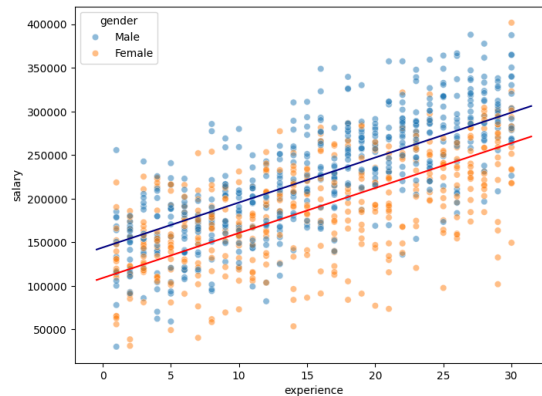
4

Question 2: Does the salary differ for men and women?

$$\text{Salary} = a + f * \text{Female} + b * \text{Experience}$$

```
# linear regression of salary on experience and female dummy
salary_experience_gender = smf.ols('salary ~ female + experience',
                                   data=salary).fit()
print(salary_experience_gender.summary())
```

		coef	std err	t	P> t	[0.025	0.975]
<i>a</i> ->	Intercept	1.441e+05	3244.879	44.401	0.000	1.38e+05	1.5e+05
<i>f</i> ->	female	-3.494e+04	2904.870	-12.027	0.000	-4.06e+04	-2.92e+04
<i>b</i> ->	experience	5154.5409	166.672	30.926	0.000	4827.473	5481.608



5

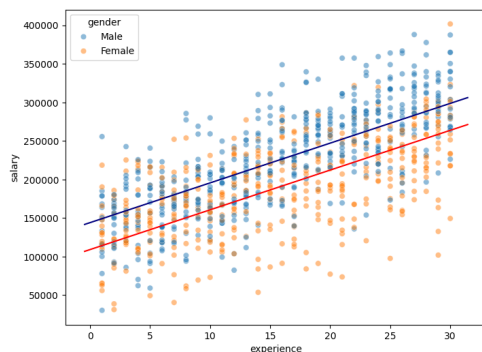
Question 2: Does the salary differ for men and women?

$$\text{Salary} = a + f * \text{Female} + b * \text{Experience}$$

		coef	std err	t	P> t	[0.025	0.975]
<i>a</i> ->	Intercept	1.441e+05	3244.879	44.401	0.000	1.38e+05	1.5e+05
<i>f</i> ->	female	-3.494e+04	2904.870	-12.027	0.000	-4.06e+04	-2.92e+04
<i>b</i> ->	experience	5154.5409	166.672	30.926	0.000	4827.473	5481.608

Men -> $\text{Salary} = a + f * 0 + b * \text{Experience} = a + b * \text{Experience}$

Women -> $\text{Salary} = a + f * 1 + b * \text{Experience} = (a + f) + b * \text{Experience}$



-34940 = amount by which
 the average salary **differs** for
 women (female=1) relative to
 men (male=0)

$(144100 - 34940) = 109,160$
 = average salary of women
 (at 0 experience)

6

Question 3: Does the effect of work experience on salary differ for men and women?

$$\text{Salary} = a + f * \text{Female} + b * \text{Experience} + c * \text{Female} * \text{Experience}$$

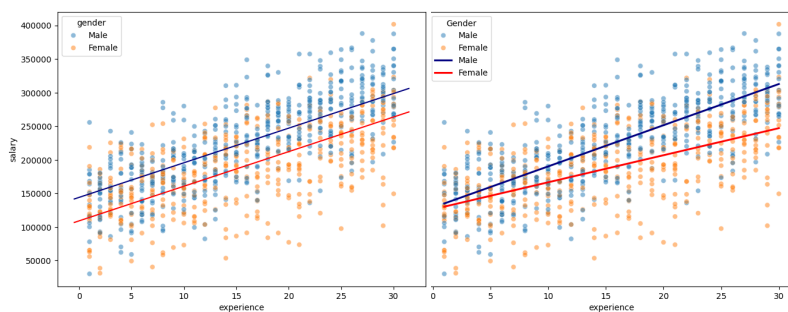
```
# Fit the linear regression model with interaction
# between 'female' and 'experience'
salary_experience_gender = smf.ols('salary ~ female * experience',
                                   data=salary).fit()
```

	coef	std err	t	P> t	[0.025	0.975]
a -> Intercept	1.289e+05	3962.822	32.531	0.000	1.21e+05	1.37e+05
f -> female	-2611.4508	5786.649	-0.451	0.652	-1.4e+04	8743.971
b -> experience	6129.6546	223.143	27.470	0.000	5691.769	6567.540
c -> female:experience	-2102.8340	327.687	-6.417	0.000	-2745.870	-1459.798

FEMALE is now statistically insignificant!

Men --> $\text{Salary} = a + f * 0 + b * \text{Experience} + c * 0 * \text{Experience} = a + b * \text{Experience}$

Women --> $\text{Salary} = a + f * 1 + b * \text{Experience} + c * 1 * \text{Experience} = (a + f) + (b + c) * \text{Experience}$



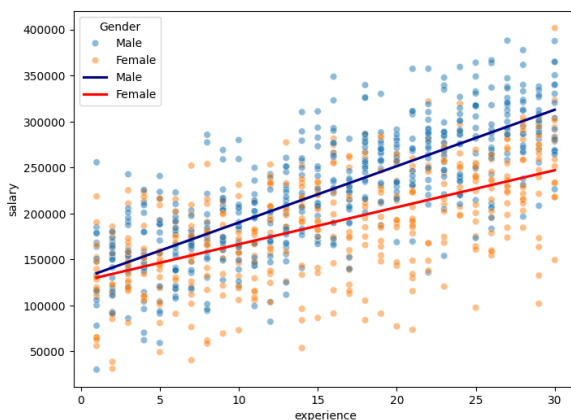
$(6130 - 2103) =$
the effect of one more year
of experience for womens' salaries

6130 = the effect of one more
year of experience for mens' salaries

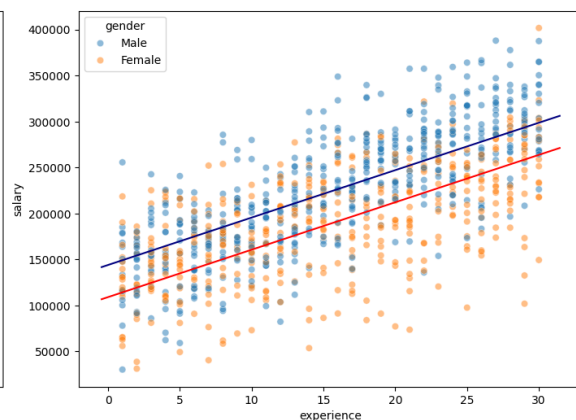
7

For analyst-driven models, our inference is only as good as our mental model is.

With Interaction



Without Interaction



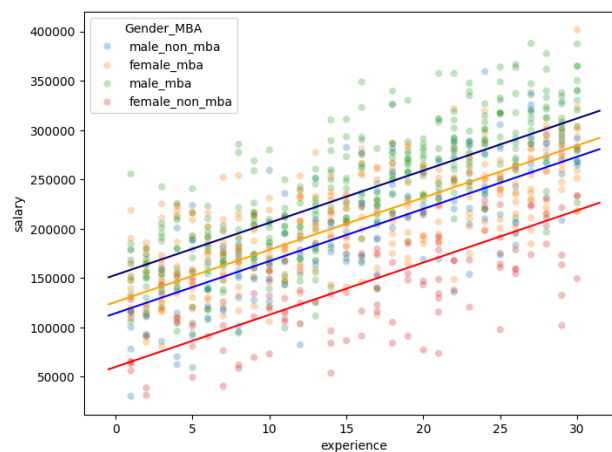
8

Interpreting Interaction Term: Two Dummy Variables

Does an MBA affect the salaries for men and women differently?

$$\text{Salary} = a + f * \text{Female} + c * \text{MBA} + d * \text{Female} * \text{MBA} + b * \text{Experience}$$

	coef	std err	t	P> t	[0.025	0.975]
<i>a</i> -> Intercept	1.144e+05	3802.692	30.082	0.000	1.07e+05	1.22e+05
<i>f</i> -> female	-5.967e+04	4661.204	-12.801	0.000	-6.88e+04	-5.05e+04
<i>c</i> -> mba	3.886e+04	3621.716	10.729	0.000	3.17e+04	4.6e+04
<i>b</i> -> female:mba	3.221e+04	5474.286	5.884	0.000	2.15e+04	4.3e+04
<i>d</i> -> experience	5289.6355	140.380	37.681	0.000	5014.160	5565.111



Does an MBA affect the salaries for men and women differently?

$$\text{Salary} = a + f * \text{Female} + c * \text{MBA} + d * \text{Female} * \text{MBA} + b * \text{Experience}$$

	coef	std err	t	P> t	[0.025	0.975]
<i>a</i> -> Intercept	1.144e+05	3802.692	30.082	0.000	1.07e+05	1.22e+05
<i>f</i> -> female	-5.967e+04	4661.204	-12.801	0.000	-6.88e+04	-5.05e+04
<i>c</i> -> mba	3.886e+04	3621.716	10.729	0.000	3.17e+04	4.6e+04
<i>d</i> -> female:mba	3.221e+04	5474.286	5.884	0.000	2.15e+04	4.3e+04
<i>b</i> -> experience	5289.6355	140.380	37.681	0.000	5014.160	5565.111

Male, no MBA --> $\text{Salary} = a + f * 0 + c * 0 + d * 0 * 0 + b * \text{Exper.} = a + b * \text{Experience}$

Male, MBA --> $\text{Salary} = a + f * 0 + c * 1 + d * 0 * 1 + b * \text{Exper.} = (a + c) + b * \text{Experience}$

Female, no MBA --> $\text{Salary} = a + f * 1 + c * 0 + d * 1 * 0 + b * \text{Exper.} = (a + f) + b * \text{Experience}$

Female, MBA --> $\text{Salary} = a + f * 1 + c * 1 + d * 1 * 1 + b * \text{Exper.} = (a + f + c + d) + b * \text{Experience}$

38860 = amount by which an

MBA affects the average

Effect of MBA on men's salaries: $(a + c) - a = c$ salary for men

Effect of MBA on women's salaries: $(a + f + c + d) - (a + f) = c + d$ 71070 = amount by which an
MBA affects the average
salary for women

11

Interpreting Multiple Interaction Terms

Multiple Dummy Variables: the Firewall Wizard Example

EXAMPLE: FIREWALL WIZARD

- Firewalls on PC are notoriously hard to manage (require knowledge of IP ports and networking)
- New “wizard” for configuring Windows firewall
- 10,000 customers are targeted with one of two Ad copies
 - Ad-copy A emphasizes “ease of use” (4,607 customers, 50% males)
 - Ad-copy B emphasizes “control/options” (5,393 customers, 50% males)
- Available data:

```
res          Is 1 if responded to offer, 0 if not
age          age of customer
numpurch     total number of purchases
totdol       total dollars spent
adB          Is 1 if Ad-copy B 'control/options', 0 if Ad-copy 'A' ease-of-use'
female       Is 1 if female, 0 if male
```

13

What predicts the response to the firewall offer?

FIREWALL EXAMPLE: RESULTS FROM A LINEAR PROBABILITY REGRESSION

```
# first regression, no interaction
firewall_no_inter = smf.ols('res ~ age + numpurch + totdol + female + adB',
                             data=firewall).fit()
print(firewall_no_inter.summary())
```

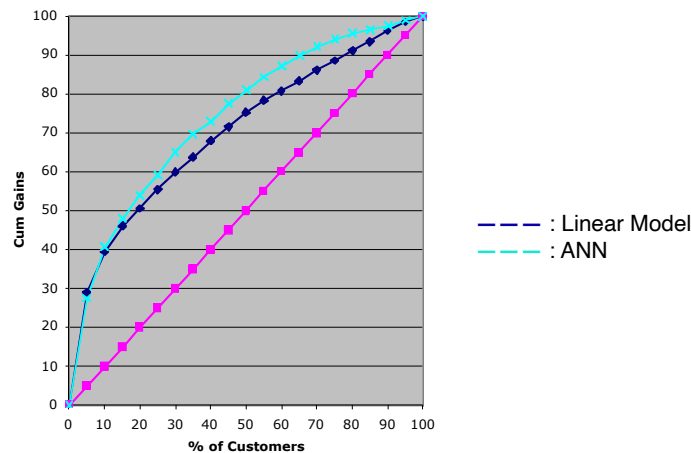
	coef	std err	t	P> t	[0.025	0.975]
Intercept	-0.1662	0.027	-6.222	0.000	-0.219	-0.114
age	0.0051	0.001	7.383	0.000	0.004	0.006
numpurch	0.0373	0.004	8.911	0.000	0.029	0.046
totdol	0.0001	7.14e-06	17.098	0.000	0.000	0.000
female	-0.0085	0.006	-1.407	0.159	-0.020	0.003
adB	-0.0033	0.006	-0.556	0.578	-0.015	0.008

- What seems to matter and what not?

14

The ANN does substantially better than the linear regression model

GAIN COMPARISON FOR ANN AND LINEAR MODEL



Are we missing something?

```
res      Is 1 if responded to offer, 0 if not
age      age of customer
numpurch total number of purchases
totdol   total dollars spent
adB      Is 1 if Ad-copy B 'control/options', 0 if Ad-copy 'A' ease-of-use'
female   Is 1 if female, 0 if male
```

15

What predicts the response to the firewall offer?

EXAMPLE 2: RESULTS FROM A NEW LINEAR REGRESSION

```
# second regression, with interaction
firewall_w_inter = smf.ols('res ~ age + numpurch + totdol + female * adB',
                           data=firewall).fit()
print(firewall_w_inter.summary())
```

	coef	std err	t	P> t	[0.025	0.975]
Intercept	-0.2056	0.027	-7.725	0.000	-0.258	-0.153
age	0.0050	0.001	7.351	0.000	0.004	0.006
numpurch	0.0375	0.004	9.036	0.000	0.029	0.046
totdol	0.0001	7.07e-06	17.326	0.000	0.000	0.000
female	0.0795	0.009	9.069	0.000	0.062	0.097
adB	0.0737	0.008	8.981	0.000	0.058	0.090
female:adB	-0.1632	0.012	-13.661	0.000	-0.187	-0.140

Interpretation:

- All customers see an ad!

- Which type of customer+ad is the baseline effect of the ad?

16

Consider the effects for different genders who receive different ads

$$a + b * \text{age} + c * \text{numpurch} + d * \text{totdol} + f * \text{female} + e * \text{adB} + g * \text{female} * \text{adB}$$

	coef	std err	t	P> t	[0.025	0.975]
Intercept	-0.2056	0.027	-7.725	0.000	-0.258	-0.153
age	0.0050	0.001	7.351	0.000	0.004	0.006
numpurch	0.0375	0.004	9.036	0.000	0.029	0.046
totdol	0.0001	7.07e-06	17.326	0.000	0.000	0.000
f -> female	0.0795	0.009	9.069	0.000	0.062	0.097
e -> adB	0.0737	0.008	8.981	0.000	0.058	0.090
g -> female:adB	-0.1632	0.012	-13.661	0.000	-0.187	-0.140

Male, adA -> $f * 0 + e * 0 + g * 0 * 0 = 0$

Male, adB -> $f * 0 + e * 1 + g * 0 * 1 = e$

Female, adA -> $f * 1 + e * 0 + g * 1 * 0 = f$

Female, adB -> $f * 1 + e * 1 + g * 1 * 1 = f + e + g$

For Men: effect of adB relative to adA $e - 0 = e \longrightarrow 0.07$

For Women: effect of adB relative to adA $(f + e + g) - f = e + g \longrightarrow 0.07 - 0.16 = -0.09$

17

What predicts the response to the firewall offer?

EXAMPLE 2: RESULTS FROM A NEW LINEAR REGRESSION

$$a + b * \text{age} + c * \text{numpurch} + d * \text{totdol} + f * \text{female} + e * \text{adB} + g * \text{female} * \text{adB}$$

	coef	std err	t	P> t	[0.025	0.975]
Intercept	-0.2056	0.027	-7.725	0.000	-0.258	-0.153
age	0.0050	0.001	7.351	0.000	0.004	0.006
numpurch	0.0375	0.004	9.036	0.000	0.029	0.046
totdol	0.0001	7.07e-06	17.326	0.000	0.000	0.000
female	0.0795	0.009	9.069	0.000	0.062	0.097
adB	0.0737	0.008	8.981	0.000	0.058	0.090
female:adB	-0.1632	0.012	-13.661	0.000	-0.187	-0.140

Interpretation:

- **For men:** the prob of purchasing **increases by 7%** when they see Ad B instead of Ad A

- **For women:** the prob of purchasing **decreases by 9%** ($=0.07-0.16$) when they see Ad B instead of Ad A

18

How to think about interaction effects (in the context of causal inference)

TIPS ON INTERACTION EFFECTS

- Don't go hunting for every possible interaction effect:
Only try interactions for which you think that there is a reason they might exist
- When you have trouble interpreting what each coefficient measures, write down the regression equation and plug in 0s and 1s

An Application of Difference-in-Differences: The Effect of Search Engine Sponsored Ads on Traffic

Search Engine Shopping Sponsored Ads

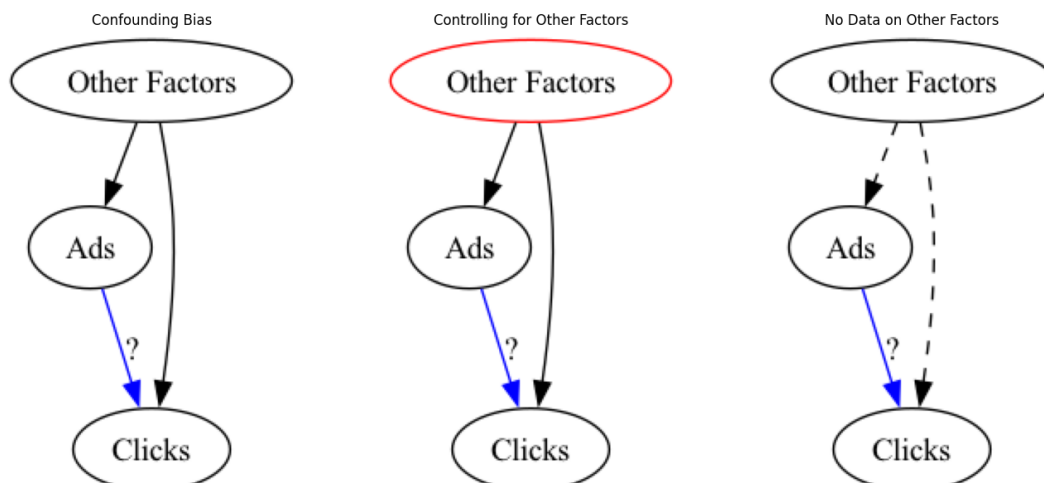
Google search results for "new balance shoes". The sponsored ads section is highlighted with a red box, showing four New Balance shoe listings with prices and shipping information. The organic search results section is highlighted with a red box, showing the New Balance website and Amazon.com listings.

Sponsored ads →

Organic search results →

Imagine an ecommerce website selling shoes and using Google/Bing for the ads. How effective are the ads?

- We have ads appearance and clicks over time and across Google/Bing, can we estimate the effect of ads on clicks?
 - Over time? Compare Google clicks between periods with/without Ads on Google
 - Between engines? Compare clicks between Google/Bing



Diff-in-Diff: Recall the coupon example

Digital coupons example

	Region 1	Region 2
Period 1	Average spending: 110	Average spending: 90
Period 2	Average spending: 120	Average spending: 140
	Target Group (A)	
	Control Group (B)	

- Difference between ? $120 - 110 = 10$
- What would happen to Region 2 if there was no promotion? $90 + 10 = 100$
- What is the effect of the promotion? $140 - 100 = \mathbf{40}$

Two Key Assumptions (out of many)

- One group's treatment status should not affect another group's outcomes (SUTVA*).
 - E.g., the treatment should not affect the control group's outcome
- Parallel trends between treatment and control groups
 - The control group can reflect how the treatment group behaves if there was no treatment.

* SUTVA: Stable Unit Treatment Value Assumption

A Natural Experiment at a Major Online Retailer

- The retailer usually advertises on both Google and Bing
- Around Labor Day in 2023 (Sep 4, 2023), a random technical glitch resulted in the retailer being unable to advertise on Google for a week (8/31/2023-9/6/2023)
- A natural experiment for diff-in-diff
 - Treatment: Glitch (not advertise), randomly happened
 - Treated and Control: #Clicks from Google vs. Bing
 - Before/During Treatment: Before the glitch and during the glitch
- Google/Bing budgets are earmarked and cannot be switched around
 - SUTVA is most likely satisfied

Data:

- One month before the glitch and 7 days during the glitch
- Daily clicks from Google and Bing

```
# load data
url = 'https://tinyurl.com/sponsored-search-ads-traffic'
sponsored_ads = pd.read_csv(url)

# randomly sample and show 10 rows of the dataframe
np.random.seed(42)
print(sponsored_ads.sample(n=10))
```

	date	glitch	clicks	log_clicks	click_source
4	2023-08-04	0	100.321000	4.608375	bing
63	2023-08-26	0	290.239990	5.670708	google
18	2023-08-18	0	96.037003	4.564733	bing
0	2023-07-31	0	103.431000	4.638905	bing
28	2023-08-28	0	93.695999	4.540055	bing
73	2023-09-06	1	293.355010	5.681384	google
10	2023-08-10	0	100.228000	4.607448	bing
34	2023-09-04	1	108.335000	4.685228	bing
12	2023-08-12	0	110.720000	4.707005	bing
55	2023-08-18	0	296.998990	5.693729	google

Over time: Compare Google clicks between periods with/without Ads on Google (before and during the glitch)

Capture weekday effects on clicks

Focus on Google, where the glitch happened

Stopping ads has no effect on #clicks?

```
# convert date to datetime
sponsored_ads['date'] = pd.to_datetime(sponsored_ads['date'])
# generate weekday variable
sponsored_ads['weekday'] = sponsored_ads['date'].dt.dayofweek

### generate dummy variable for clicks from google
sponsored_ads['google'] = np.where(sponsored_ads['click_source'] == 'google', 1, 0)

### generate interaction term between google and glitch, the treatment variable
sponsored_ads['google_glitch'] = sponsored_ads['google'] * sponsored_ads['glitch']
```

```
### 1. simple regression of glitch using google data only
#### This regression tells the change of google-clicks before and during the glitch
time_diff = \
    smf.ols('log_clicks ~ glitch + C(weekday)', \
            data=sponsored_ads[sponsored_ads['google']==1]).fit()
```

	coef	std err	t	P> t	[0.025	0.975]
Intercept	5.7148	0.007	806.349	0.000	5.700	5.729
C(weekday) [T.1]	-0.0167	0.010	-1.691	0.101	-0.037	0.003
C(weekday) [T.2]	-0.0082	0.010	-0.791	0.435	-0.029	0.013
C(weekday) [T.3]	-0.0272	0.010	-2.626	0.014	-0.048	-0.006
C(weekday) [T.4]	-0.0054	0.010	-0.517	0.609	-0.027	0.016
C(weekday) [T.5]	-0.0247	0.010	-2.380	0.024	-0.046	-0.003
C(weekday) [T.6]	-0.0151	0.010	-1.456	0.156	-0.036	0.006
glitch	-0.0112	0.007	-1.562	0.129	-0.026	0.003

Between engines: Compare clicks between Google/Bing

Focus on the week with the glitch

Stopping ads actually increase #clicks?

```
### 2. simple regression of glitch using google data only
#### This regression tells the comparison of clicks between google and bing during the glitch
source_diff = \
    smf.ols('log_clicks ~ google + C(weekday)', \
            data=sponsored_ads[sponsored_ads['glitch']==1]).fit()
```

	coef	std err	t	P> t	[0.025	0.975]
Intercept	4.6706	0.033	141.990	0.000	4.590	4.751
C(weekday) [T.1]	-0.0519	0.044	-1.192	0.278	-0.158	0.055
C(weekday) [T.2]	-0.0323	0.044	-0.743	0.485	-0.139	0.074
C(weekday) [T.3]	0.0070	0.044	0.160	0.878	-0.099	0.113
C(weekday) [T.4]	-0.0143	0.044	-0.329	0.754	-0.121	0.092
C(weekday) [T.5]	-0.0801	0.044	-1.841	0.115	-0.187	0.026
C(weekday) [T.6]	-0.0542	0.044	-1.246	0.259	-0.161	0.052
google	1.0514	0.023	45.203	0.000	0.994	1.108

Diff-in-Diff

```
### generate interaction term between google and glitch, the treatment variable
sponsored_ads['google_glitch'] = sponsored_ads['google'] * sponsored_ads['glitch']
```

Using data of both engines and both periods

```
### 3. run the diff-in-diff regression
diff_in_diff = \
    smf.ols('log_clicks ~ google + glitch + google_glitch + C(weekday)', \
            data=sponsored_ads).fit()

### print the summary of the regressions
print(diff_in_diff.summary())
```

	coef	std err	t	P> t	[0.025	0.975]
Intercept	4.6084	0.013	366.394	0.000	4.583	4.634
C(weekday) [T.1]	-0.0219	0.016	-1.363	0.178	-0.054	0.010
C(weekday) [T.2]	-0.0084	0.017	-0.500	0.619	-0.042	0.025
C(weekday) [T.3]	0.0035	0.017	0.207	0.836	-0.030	0.037
C(weekday) [T.4]	-0.0082	0.017	-0.487	0.628	-0.042	0.025
C(weekday) [T.5]	-0.0076	0.017	-0.454	0.651	-0.041	0.026
C(weekday) [T.6]	-0.0061	0.017	-0.361	0.720	-0.040	0.028
google	1.1002	0.010	108.432	0.000	1.080	1.120
glitch	0.0369	0.017	2.236	0.029	0.004	0.070
google_glitch	-0.0488	0.023	-2.090	0.041	-0.095	-0.002

Interpretation?

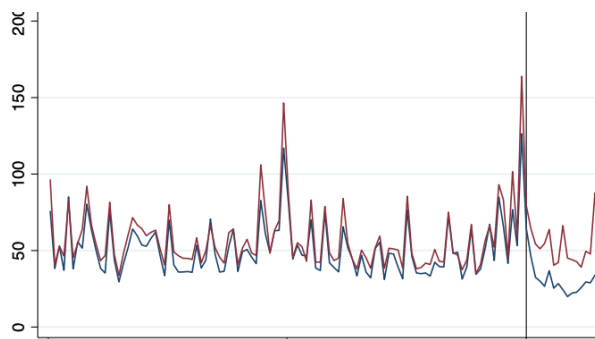
Parallel Trends: The treated and the control groups have similar variations before Treatment

- Why do we need the parallel trends?

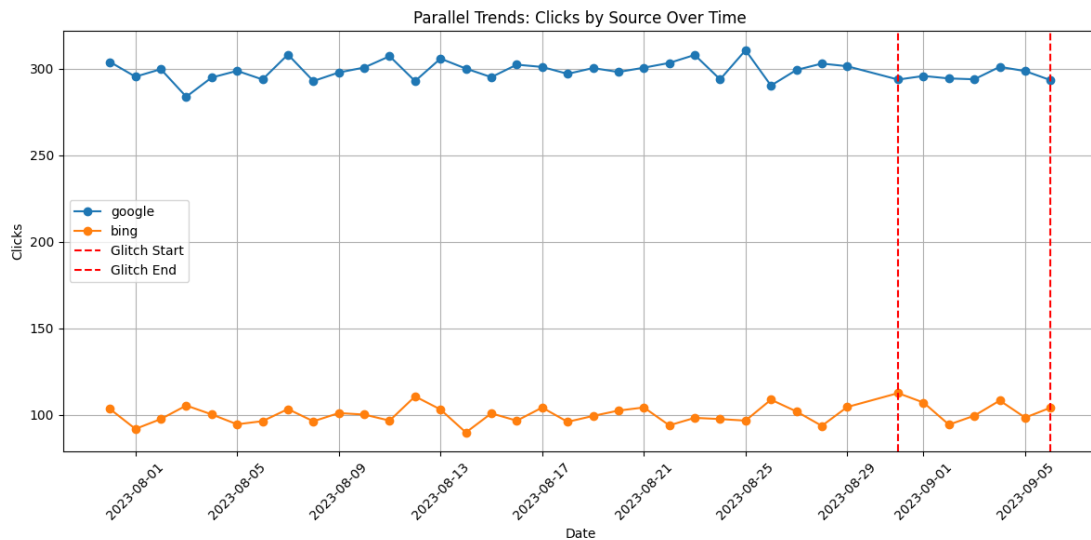
$$E(Y|T=1) - E(Y|T=0) = E(Y_1 - Y_0|T=1) \quad \leftarrow \text{ATT}$$

$$+ E(Y_0|T=1) - E(Y_0|T=0) \quad \leftarrow \text{Bias}$$

We need the control group to show us what would happen if no treatment!



However, the graph sometimes can be messy



Parallel Trend? More systematic approaches: Time Interaction

- Interact the treated group with time trends or time dummies before the treatment
 - If the interaction terms are insignificant, it implies the treated group has no significant difference from the control group

A time index,
0, 1, 2, 3, ...

Before the
Glitch

```
## generate a time trend variable using the date variable
sponsored_ads['time_trend'] = (sponsored_ads['date'] - sponsored_ads['date'].min()).dt.days

parallel_trend1 = \
    smf.ols('log_clicks ~ google * time_trend', \
            data=sponsored_ads[sponsored_ads['glitch']==0]).fit()
print(parallel_trend1.summary())
```

	coef	std err	t	P> t	[0.025	0.975]
Intercept	4.5944	0.013	344.333	0.000	4.568	4.621
google	1.0989	0.019	58.238	0.000	1.061	1.137
time_trend	0.0005	0.001	0.586	0.560	-0.001	0.002
google:time_trend	8.354e-05	0.001	0.075	0.941	-0.002	0.002

Interpretation?

Another Standard Robustness Check: Placebo Test

- Basic idea: If we artificially move the treatment to a pre-treatment period, *it should have no effect*. Otherwise, it indicates non-parallel trends or other confounding factors the diff-in-diff cannot address

Define the placebo treatment

```
### first generate the placebo variable--moving the glitch to 7/31/2023 to 8/6/2023
sponsored_ads['placebo_glitch'] = \
    np.where((sponsored_ads['date'] >= '2023-07-31') & (sponsored_ads['date'] <= '2023-08-06'), 1, 0)

### generate the interaction term between google and placebo_glitch
sponsored_ads['google_placebo_glitch'] = sponsored_ads['google'] * sponsored_ads['placebo_glitch']

### run the diff-in-diff regression
diff_in_diff_placebo = \
    smf.ols('log_clicks ~ google + placebo_glitch + google_placebo_glitch + C(weekday)', \
            data=sponsored_ads).fit()

# # # print the summary of the regressions
print(diff_in_diff_placebo.summary())
```

	coef	std err	t	P> t	[0.025	0.975]
google	1.0889	0.010	104.935	0.000	1.068	1.110
placebo_glitch	-0.0234	0.017	-1.385	0.171	-0.057	0.010
google_placebo_glitch	0.0109	0.024	0.457	0.649	-0.037	0.059

Interpretation?