

ROI-on-Sponsored-Search

Jhan-Syuan Lin

March 19, 2023

Business Setup

Bazaar.com is the top online store in the United States, using sponsored search advertisements on multiple platforms. It publishes advertising in response to keywords entered by online customers and divides them into branded and nonbranded categories. ‘Bazaar shoes’ and ‘Bazaar guitar’ are examples of brand keywords that include the brand name. Nonbranded keywords include generic terms such as ‘shoes’ and ‘guitar.’

Using traffic data from several platforms, bazaar’s marketing team calculated a 320% ROI for sponsored search advertisements. This result is troublesome since visitors who searched for ‘Bazaar’ already intended to visit Bazaar.com; therefore, we question the usefulness of branded keyword advertisements. To achieve our aim of understanding the causal inference of search advertising and its efficacy, the following analysis will be performed: * What’s wrong with Bob’s ROI analysis? * Define the Treatment and Control. * Consider a First Difference Estimate. * Calculate the Difference-in-Differences. * Given the Treatment Effect Estimate, Fix Bob’s RoI Calculation.

```
library(dplyr)
library(ggplot2)
library(plm)

# Set working directory to source file location
setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
```

Questions

(a) What is Wrong with Bob's RoI Calculation?

As case mentioned, the 12% conversion rate we observed is not purely based on sponsored traffic but on both the sponsored and organic links. Therefore, we must isolate the conversion rate for sponsored ads only to calculate the right ROI. Given the sponsored ads on branded keywords, people who would have used organic search could also use sponsored ads to reach the website. These people usually have a higher conversion rate since they are already familiar with the brand. This fact could lead to a wrong conclusion about the conversion rate in sponsored ads.

Besides, the margin per conversion is \$21. This number is also biased since it is calculated by a combination of both sponsored and organic links. The actual margin could even be lower for those who click on the sponsored ads since they probably are still in the awareness phase.

(b) Define the Treatment and Control. What is the unit of observation here?

Define the treatment. Which unit(s) are treated and which is / are control?

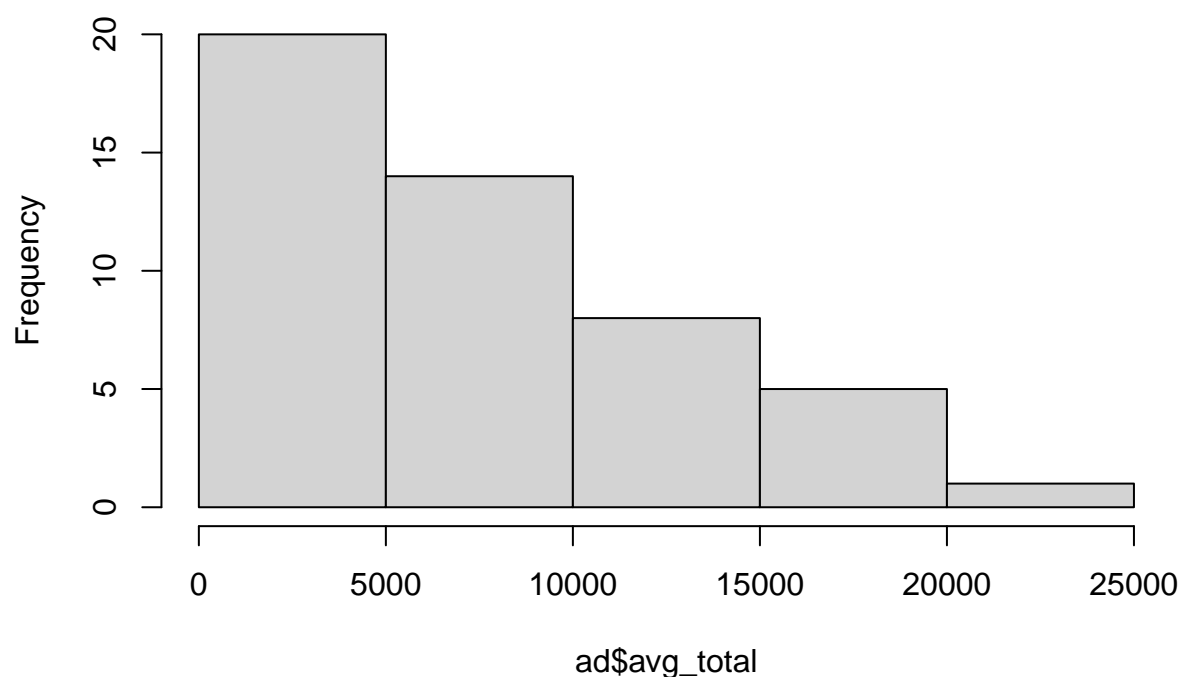
- Unit of observation: Weekly average clicks number on each platform.
- Treatment: Suspend sponsored ad campaign (Ans: Google platform data after week 9)
- Treatment Group: Average clicks number of Google platform
- Control Group: Average clicks number of other platforms

(c) Consider a First Difference Estimate.

```
ad = read.csv("did_sponsored_ads.csv")
ad$avg_total = ad$avg_org + ad$avg_spons

hist(ad$avg_total) # Although the y look pretty skewed
```

Histogram of ad\$avg_total



```
# Create Dummy Variables
ad = ad%>%mutate(after = ifelse(week<10, 0, 1))
ad = ad%>%mutate(treatment = ifelse(id==3, 1, 0))

# Create treatment subset
google = ad%>%filter(id==3)

# Calculate the mean avg_total in the two time periods(after)
google %>%
  group_by(after)%>%
  summarise(avg_week_total = mean(avg_total),
            avg_week_spons = mean(avg_spons),
            avg_week_org = mean(avg_org))
```

```
## # A tibble: 2 x 4
```

```
##   after avg_week_total avg_week_spons avg_week_org
```

```
##      <dbl>          <dbl>          <dbl>          <dbl>
## 1      0            8390.            6123.            2267.
## 2      1            6544              0            6544
```

```
# First difference
```

```
fd_model = lm(avg_total ~ after, data=google)
summary(fd_model)
```

```
##
## Call:
## lm(formula = avg_total ~ after, data = google)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7003.9 -2630.1  -172.5   2088.4   8625.1
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      8390         1598   5.252 0.000373 ***
## after           -1846         3195  -0.578 0.576238
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4793 on 10 degrees of freedom
## Multiple R-squared:  0.0323, Adjusted R-squared:  -0.06447
## F-statistic: 0.3337 on 1 and 10 DF,  p-value: 0.5762
```

Although the histogram of total website traffic looks pretty skewed, we will stick with the non transformation model to simplify the explanation. With the first difference method, we can see that the treatment effect (no sponsored ad) causes around -1846 decrease in total web traffic for the after period of the Google platform. The % change of clicks due to the absence of sponsored ads $(6544-8390) / 8390$, which is around 22%. However, we must interpret this result with caution because the p-value is greater than 0.05, indicating no evidence that this treatment affects the average total click.

The reason why this number is not solely reliable is that we ignore the natural variant of the website traffic. That said, perhaps in the post-period, the website traffic shows a significantly different trend compared to the pre-period. The estimation with this model could not capture this element and hence might lead to a wrong conclusion.

(d) Calculate the Difference-in-Differences.

```
# Check Parallel Trend
```

```
summary(lm(avg_total ~ factor(week)*treatment, data=ad))
```

```
##
```

```
## Call:
```

```
## lm(formula = avg_total ~ factor(week) * treatment, data = ad)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -8710.7  -111.8    87.3   1422.3  6586.3
```

```
##
```

```
## Coefficients:
```

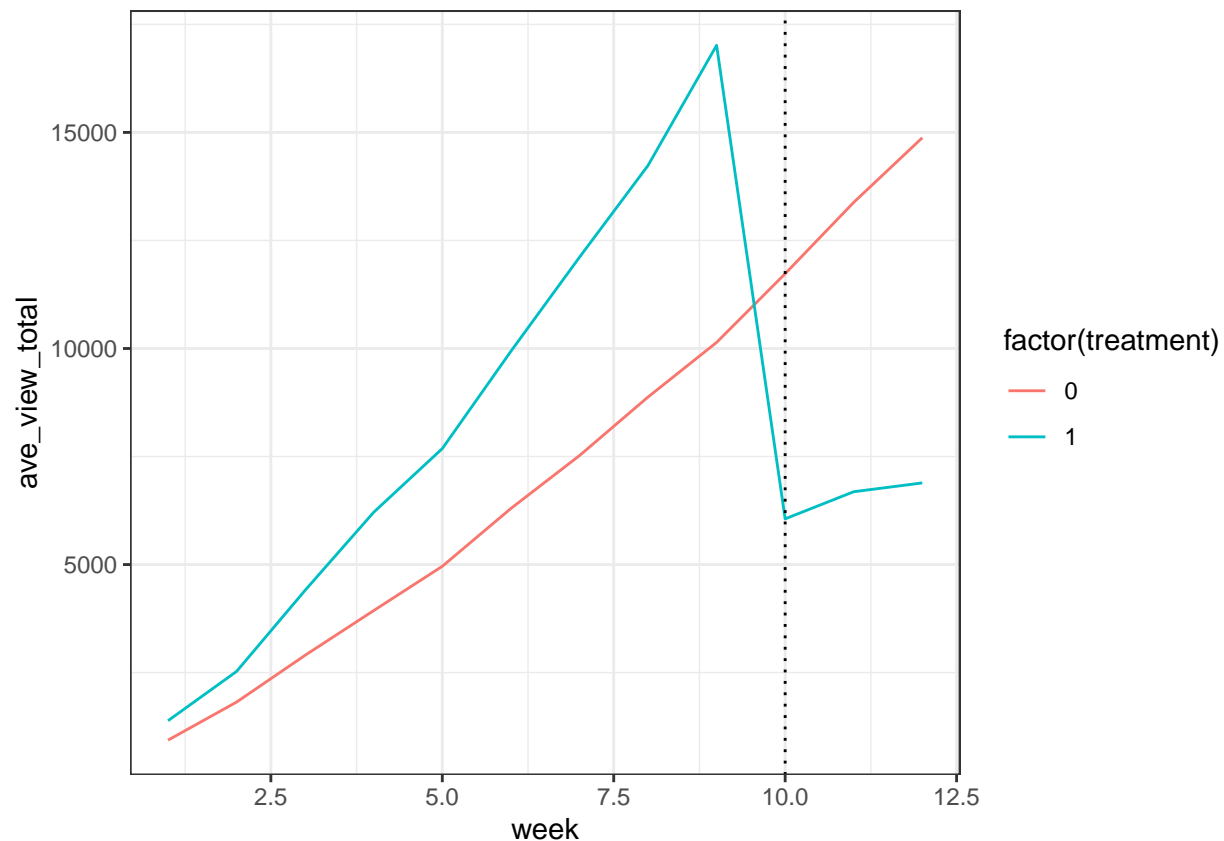
```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)          936.3     2465.2   0.380 0.707414
## factor(week)2          881.3     3486.3   0.253 0.802574
## factor(week)3         1964.7     3486.3   0.564 0.578291
## factor(week)4         2998.3     3486.3   0.860 0.398274
## factor(week)5         4023.3     3486.3   1.154 0.259840
## factor(week)6         5361.0     3486.3   1.538 0.137190
## factor(week)7         6584.7     3486.3   1.889 0.071069 .
## factor(week)8         7940.0     3486.3   2.278 0.031955 *
## factor(week)9         9204.3     3486.3   2.640 0.014337 *
## factor(week)10        10794.3     3486.3   3.096 0.004932 **
## factor(week)11        12445.3     3486.3   3.570 0.001550 **
## factor(week)12        13940.3     3486.3   3.999 0.000529 ***
## treatment              449.7     4930.3   0.091 0.928087
```

```
## factor(week)2:treatment      259.7      6972.5      0.037 0.970600
## factor(week)3:treatment      1055.3      6972.5      0.151 0.880960
## factor(week)4:treatment      1826.7      6972.5      0.262 0.795571
## factor(week)5:treatment      2274.7      6972.5      0.326 0.747075
## factor(week)6:treatment      3187.0      6972.5      0.457 0.651723
## factor(week)7:treatment      4140.3      6972.5      0.594 0.558196
## factor(week)8:treatment      4909.0      6972.5      0.704 0.488177
## factor(week)9:treatment      6424.7      6972.5      0.921 0.365997
## factor(week)10:treatment     -6122.3      6972.5     -0.878 0.388613
## factor(week)11:treatment     -7146.3      6972.5     -1.025 0.315616
## factor(week)12:treatment     -8437.3      6972.5     -1.210 0.238030
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4270 on 24 degrees of freedom
## Multiple R-squared:  0.6819, Adjusted R-squared:  0.3771
## F-statistic: 2.237 on 23 and 24 DF,  p-value: 0.0278
```

```
# Group data by week and treatment and calculate average values for plotting
week_ave = ad %>% group_by(week, treatment) %>% summarise(ave_view_total = mean(avg_total),
                                                         ave_view_org = mean(avg_org),
                                                         ave_view_spons = mean(avg_spons))
```

```
## 'summarise()' has grouped output by 'week'. You can override using the
## '.groups' argument.
```

```
ggplot(week_ave, aes(x = week, y = ave_view_total, color = factor(treatment))) +
  geom_line() +
  geom_vline(xintercept = 10, linetype='dotted') +
  theme_bw()
```



```
# Calculate the mean avg_total in the two time periods(after)
```

```
ad %>%
```

```
  group_by(treatment, after)%>%
```

```
  summarise(avg_week_total = mean(avg_total),
```

```
            avg_week_spons = mean(avg_spons),
```

```
            avg_week_org = mean(avg_org))
```

```
## 'summarise()' has grouped output by 'treatment'. You can override using the
```

```
## '.groups' argument.
```

```
## # A tibble: 4 x 5
```

```
## # Groups:   treatment [2]
```

```
##   treatment after avg_week_total avg_week_spons avg_week_org
```

```
##      <dbl> <dbl>          <dbl>          <dbl>          <dbl>
```

```
## 1      0    0          5265.          3775.          1490.
```

```
## 2      0    1          13330.         9856.          3474.
```

## 3	1	0	8390.	6123.	2267.
## 4	1	1	6544	0	6544

```
# Did Model
```

```
did_model = plm(avg_total ~ treatment*after,
                 data=ad,
                 model='within',
                 effect='twoways',
                 index=c('id','week'))
summary(did_model)
```

```
## Twoways effects Within Model
```

```
##
```

```
## Call:
```

```
## plm(formula = avg_total ~ treatment * after, data = ad, effect = "twoways",
```

```
##      model = "within", index = c("id", "week"))
```

```
##
```

```
## Balanced Panel: n = 4, T = 12, N = 48
```

```
##
```

```
## Residuals:
```

```
##      Min.   1st Qu.   Median   3rd Qu.    Max.
```

```
## -4415.19 -963.23 -242.25  900.76 4216.45
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t-value Pr(>|t|)
```

```
## treatment:after -9910.6      1728.3 -5.7344 2.345e-06 ***
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
```

```
## Total Sum of Squares:    327040000
```

```
## Residual Sum of Squares: 161290000
```

```
## R-Squared:      0.50681
```

```
## Adj. R-Squared: 0.27562
```


F-statistic: 32.8834 on 1 and 32 DF, p-value: 2.3449e-06

```
# The real % Loss of clicks due to absence of sponsored ads with DiD
# (6544-8390) / 8390 - (13330-5265) / 5265 # 175% decrease
```

Looking at the graph and the interaction terms of the dynamic DiD model, we can see that the parallel trend assumption does not hold with the data. However, we will proceed with our analysis. With the DiD model, we can discover that the difference in difference effect of the treatment is -9910.6, which is way lower than the coefficient we estimate by the First Difference method. This shows the real impact of suspending sponsored ads on branded keywords. More specifically, this DiD model captures the difference in total traffic for the post-period with and without the treatment effect, which the first difference model could not capture.

(e) New RoI calculation.

This ROI calculation is still based on the information provided by Bob (e.g the conversion rate and the margin), which might be not very accurate as we discussed in question (a).

- Incremental weekly traffic attribute to sponsored ad: 9911
- Incremental gain from these clicks: $9911 * 0.12 * 21 = 24975.72$
- Average weekly clicks from sponsored search: 6123
- Weekly cost of sponsored search: $6123 * 0.60 = 3673.8$

$$ROI = (24975.72 - 3673.8) / 3673.8 = 580\%$$

Rev: $9910.59 * 2.25 = 24974$ Cost: Bazzar paid for 12681 sponsored ads in Week 9, $12681 * 0.6 = 7608$ (call it 8000) ROI: $(24974 - 8000) / 8000 \approx 215\%$