

# Causal Impact Estimation of Sponsored Ads for Bazaar.com

## Contents

Business Overview . . . . .	1
Experiment . . . . .	1
Executive summary . . . . .	1
Threats to Causal Inference . . . . .	2
Data Exploration and Overview . . . . .	2
Data Analysis . . . . .	3

## Business Overview

Bazaar.com is a leading online retailer with strong online advertising presence through paid search ads on major search engines such as Google, Bing, Yahoo and Ask. Paid ads are broadly classified into two different buckets based on the keyword bid - Branded and Non-Branded. Branded key words are ones where users search for products with part of their query containing the word 'Bazzar'. The rest are classified as non-branded keywords. Bazaar.com's marketing analytics team claim that their ad spending results in 320% ROI. However, with further discussions, they suspect the correctness of RoI calculation and want to deep-dive and analyse the actual RoI. Given the fact an user already intends to visit Bazaar.com by searching with the word 'Bazaar', they are also curious about the effectiveness of branded keyword ads and why should they invest in branded keywords.

As Analytics Consultants, we are now at task to help the marketing team to understand and answer the following key questions: 1. What's wrong with current ROI analysis? 2. What is the causal effectiveness of branded keyword ads in driving traffic to Bazaar.com? Should they continue investments? 3. Given the effectiveness of branded keyword ads, what's the corrected RoI?

## Experiment

Bazzar.com faced a technical glitch in Google search engine for a brief period of time and that gave a good experimental setup to analyse and estimate the effectiveness of sponsored search ads. We extracted 12 weeks data on number of visits to Bazzar.com through sponsored branded keyword ads as well as through organic search results. The technical glitch occured on week 10, leading to no sponsored ad starting 10th week uptill 12th week. On the other hand, Bing, Yahoo and Ask ran uninterrupted sponsored ads for all 12 weeks. Given the ad strategies, keyword bidding and mix of potential customers visting the website are same across Google and other search engine, it would be safe to assume this as an experimental setup to detect the causal impact of the ads. Treatment here is discontinuing of sponsored search ads in Google during weeks 10,11,12. Google acts as the treatment group and other search engines such as Bing, Yahoo and Ask act as the control group.

## Executive summary

Investigating the current RoI calculations, we found that the calculation is not valid for two reasons - Inflated Revenue and Opportunity cost involved.

**Inflated revenue argument** Among all customers who clicked sponsored ads, only a portion of them would have visited Bazaar's website due to sponsored ad. The rest would still visit the website through organic links in the absence of sponsored ads, as they have already showed intent by typing the branded keyword as part of the search. Therefore, the marketing team should not account all the clicks from users coming via sponsored ads to be truly causal, which otherwise would lead to overestimated revenue and ROI. \* We found that only ~73% of traffic through sponsored ads were truly a result of displaying a sponsored ad

**Opportunity Cost** Bazaar.com pays \$0.6 for the rest ~27% of visitors, who would have anyway landed on the site in the absence of an ad. This is an opportunity cost that Bazaar.com could have used for other marketing investments.

To estimate the appropriate causal effect of sponsored ad on Bazaar.com's traffic, we pursued difference-in-difference method of estimation, which proceeds through the following two steps: 1) Calculate first difference for average weekly total traffic (Ads + Organic) between the timeperiods before and after the technical glitch. This determines the raw effect of sponsored ads within google but not the overall treatment effect 2) Compare and difference the first level pre-post difference in Google and the first level pre-post difference in other search engines. This determines the true incremental effect as this step handles possible confounders like seasonal variations across weeks as well as market factors.

From the results, we found that Bazaar.com would loose ~23K clicks per week on an average, in the absence of running sponsored branded ads. This is ~73% of weekly traffic from sponsored ad, leading us to conclude that the rest ~27% traffic visits Bazaar.com through organic results. On an overall, Bazaar.com would see ~15% drop in its overall website traffic if it stops investing on sponsored branded ads.

Taking this into account, the corrected RoI from sponsored branded ads is determined to be ~205%. Though the adjusted ROI is lower than earlier estimate of 320%, it is still substantially high return for the amount invested. Hence, it is only appropriate to continue advertisting through branded search ads.

## Threats to Causal Inference

To establish confidence in the estimation, we also discuss the common threats to causal inference.

**Selection Bias:** The primary threat to selection bias is observed while selecting the treatment and control group. In this experiment, Bazaar.com follows similar ad strategies, genre of keyword bidding and also has the similar mix of potential customers across Bing and Google. We also assume that similar set up is followed for Yahoo and Ask. So having Google search engine as the treatment group does not cause any selection bias.

**Omitted Variable Bias:** Current dataset covers only traffic/ impressions related information as part of this analysis. There might be few other external variables that compromises the treatement itself and subsequently the overall impressions count. Eg. Some major power outage forcing people not having access to internet.

**Simulaneity Bias:** There is no simultaneity bias in the experiment since we do not find any reasons to believe that visting the website would make a user to click more of search ads.

**Measurement Error:** We assume that there is no measurement error as the only variable captured at an user level being traffic to the website would not be difficult to track.

**Awarness:** We assume that consumers would not be aware of being part of experiment as users tend to use mainly their default browser which would mostly ensure clear google users from bing users.

## Data Exploration and Overview

The dataset used for the analysis is average weekly traffic data from sponsored branded ads through all the four search engines for a duration of 12 weeks. It captures the overall traffic and also contains the split for sponsored and organic traffic.

## Importing necessary packages

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

## Data Transformation

```
data = read.csv('sponsored_ads.csv')
```

After loading the data, we created an “after” flag that indicates 1 for treatment duration and 0 for pre treatment time frame. Alongside, we also created a “treatment” flag for test vs. control identification.

```
treatment_week = c(10,11,12)
data1 <- data %>% mutate(treatment=ifelse(search_engine=='goog',1,0),
                        after=ifelse(week %in% treatment_week,1,0))
```

## Data Analysis

### 1. What’s wrong with current ROI analysis?

Problem with current approach in calculating the RoI can be segmented into two different buckets. The first argument is with respect to the causality assumption of all clicks coming from sponsored ads and the other argument is from the perspective of opportunity cost. Cases for both these arguments are presented below

**Inflated revenue argument** Among all customers who clicked sponsored ads, only a portion of them would have visited Bazaar’s website due to sponsored ad. The rest would still visit the website through organic links in the absence of sponsored ads, as they have already showed intent by typing the branded keyword as part of the search. Therefore, the marketing team should not account all the clicks from users coming via sponsored ads to be truly causal, which otherwise would lead to overestimated revenue and ROI.

Below is a quick example to demonstrate the use case 1). If 100 users reach Bazaar.com through sponsored ads, the tea, would have estimated the ROI based on for all 100 users whose probability of purchase is 0.12 and avg. margin per customer is \$21  $ROI = (\$21 * 0.12 * 100) - (100 * \$0.6) / (100 * \$0.6) = 251\%$

```
ROI = (21 * 0.12 * 100) - (100*0.6) / (100 * 0.6)
ROI
```

```
## [1] 251
```

Whereas in reality only 20(say) would have been truly caused due to the ads. Adj.  $ROI = (21 * 0.12 * 20) - (100 * 0.6) / (100 * 0.6) = 49.4\%$

```
AdjROI = (21 * 0.12 * 20) - (100*0.6) / (100 * 0.6)
AdjROI
```

```
## [1] 49.4
```

**Opportunity Cost** Taking cue from above example, Bazaar paid \$0.6 for around 80 customers who would have anyway landed in the site. This  $(80 * \$0.6)$  is the opportunity cost that Bazaar.com could have used this amount for some other marketing tactics.

## 2. What is the causal effectiveness of branded keyword ads in driving traffic to Bazaar.com? Should they continue investments?

The dataset contains information at week level for four different search engines namely Google, Bing, Yahoo and Ask. The unit of observation is at week level. Treatment here is technical glitch leading to discontinuation of sponsored search ads in Google during weeks 10,11,12. Week 10,11,12 in google search engine are the treated units whereas Bing and other search engines act as the control group

### Trying a simple pre-post estimator

```
# filter data for google
data_google<- data1 %>%
  filter(search_engine=="goog")

l1 <- lm(log(tot_traffic) ~ after,data=data_google)
summary(l1)

##
## Call:
## lm(formula = log(tot_traffic) ~ after, data = data_google)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.024436 -0.011497 -0.003054  0.014618  0.022675
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  11.961524   0.005975  2001.99 < 2e-16 ***
## after        -0.161791   0.011950  -13.54 9.32e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.01792 on 10 degrees of freedom
## Multiple R-squared:  0.9483, Adjusted R-squared:  0.9431
## F-statistic: 183.3 on 1 and 10 DF,  p-value: 9.318e-08
```

**Interpretation:** With no sponsored search ads in google, we observe a decline of 16.17% in weekly total traffic to the website

\*\*This before-after estimation does not give the true treatment effect because: + 1. Before and after estimation for only the treated dataset only tells us the effect of sponsored ads within google but not the overall treatment effect + 2. It also assumes all other external factors stay constant and the treatment is the only varying event. This assumption would be violated if there are any changes in the external factors such as emergence of new competitors +3. It also assumes that week 1-9 are not fundamentally different from week 10-12 in the sense there are no seasonal variation across weeks

With all these assumptions, it would be difficult to establish causality of the sponsored ads in google. Hence we resort to Difference in Difference(DiD).

### Difference-in-Difference estimation

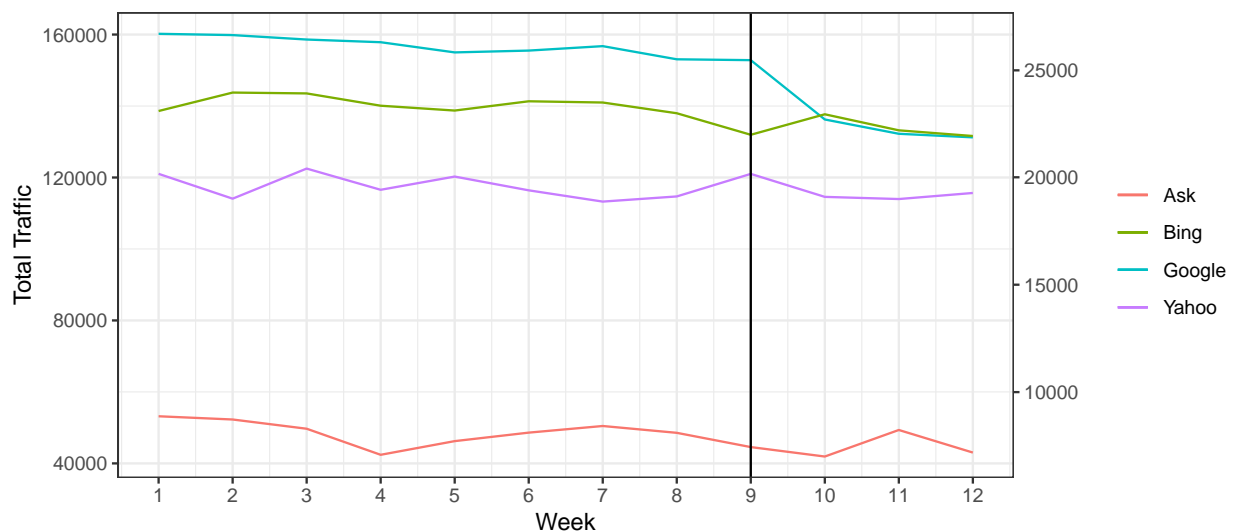
## Checking for Parallel Trend

To answer the above question, we would be using difference in difference method to estimate causal impact of the search ads. Before using the method, we check the parallel assumption trend ie; treatment and control group should have similar behaviour in the pre time frame. We first start by plotting the trend line and visualising the charts.

```
temp1 = data1 %>% filter(search_engine %in% c('bing')) %>% select(week,tot_traffic)
temp2= data1 %>% filter(search_engine %in% c('yahoo')) %>% select(week,tot_traffic)
temp3 = data1 %>% filter(search_engine %in% c('ask')) %>% select(week,tot_traffic)
ggplot(data1 %>% filter(search_engine == 'goog'),aes(x=week, y= tot_traffic, color = 'Google')) +
  geom_line() +
  geom_line(aes(x=week, y= tot_traffic*6, color = 'Bing'),data = temp1) +
  geom_line(aes(x=week, y= tot_traffic*6, color = 'Yahoo'),data = temp2) +
  geom_line(aes(x=week, y= tot_traffic*6, color = 'Ask'),data = temp3) +
  geom_vline(xintercept = 9,type = '-')+
  xlim(1,12)+
  scale_y_continuous(sec.axis = sec_axis(~./6))+
  scale_x_continuous(breaks = seq(1, 12, by = 1))+
  labs(y="Total Traffic", x = "Week")+
  theme_bw()+
  theme(legend.title=element_blank())
```

```
## Warning: Ignoring unknown parameters: type
```

```
## Scale for 'x' is already present. Adding another scale for 'x', which
## will replace the existing scale.
```



From the above visualization, it is difficult to establish the parallel trend between the search engines. Hence, we resort to dynamic DiD to check the pre trend parallel assumption

```
did_dyn1 <- lm(tot_traffic ~ treatment + factor(week) + treatment * factor(week),data=data1)
summary(did_dyn1)
```

```
##
```

```
## Call:
## lm(formula = tot_traffic ~ treatment + factor(week) + treatment *
##     factor(week), data = data1)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9546.2 -2062.1   891.9  4095.1  6733.3
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      17378.6      4578.0   3.796 0.000881 ***
## treatment       142824.5     9155.9  15.599 4.61e-14 ***
## factor(week)2      -149.8      6474.2  -0.023 0.981737
## factor(week)3       161.4      6474.2   0.025 0.980320
## factor(week)4      -764.5      6474.2  -0.118 0.906987
## factor(week)5      -421.8      6474.2  -0.065 0.948589
## factor(week)6      -359.7      6474.2  -0.056 0.956153
## factor(week)7      -451.1      6474.2  -0.070 0.945031
## factor(week)8      -643.4      6474.2  -0.099 0.921665
## factor(week)9      -849.5      6474.2  -0.131 0.896706
## factor(week)10     -1035.1     6474.2  -0.160 0.874312
## factor(week)11      -905.4      6474.2  -0.140 0.889950
## factor(week)12     -1250.6     6474.2  -0.193 0.848459
## treatment:factor(week)2    -203.0     12948.5  -0.016 0.987622
## treatment:factor(week)3   -1756.5     12948.5  -0.136 0.893227
## treatment:factor(week)4   -1586.4     12948.5  -0.123 0.903511
## treatment:factor(week)5   -4779.4     12948.5  -0.369 0.715283
## treatment:factor(week)6  -4331.7     12948.5  -0.335 0.740881
## treatment:factor(week)7  -2993.4     12948.5  -0.231 0.819135
## treatment:factor(week)8  -6483.1     12948.5  -0.501 0.621154
## treatment:factor(week)9  -6522.7     12948.5  -0.504 0.619032
## treatment:factor(week)10 -22933.8     12948.5  -1.771 0.089235 .
## treatment:factor(week)11 -27063.8     12948.5  -2.090 0.047378 *
## treatment:factor(week)12 -27718.6     12948.5  -2.141 0.042664 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7929 on 24 degrees of freedom
## Multiple R-squared:  0.9908, Adjusted R-squared:  0.982
## F-statistic: 112.6 on 23 and 24 DF, p-value: < 2.2e-16
```

From the above results, it is evident that starting week 10, the p-value is significant. It implies, prior to week 10 there is no significant behaviour in the trend lines across search engines and once the treatment starts they are significantly different.

With the parallel trend assumption established, we now perform a Difference in Difference regression between treatment and control groups to estimate the true causality of the sponsored ads. Independent variables for the DiD would be + Treatment flag( indicating test vs. control) + after flag (indicating pre vs. post) and + Interaction term between Treatment flag and the after flag.

```
# test treatment effect using DiD, this is the real loss in clicks due to sponsored ads
did <- lm(tot_traffic ~ treatment + after + treatment * after, data=data1)
summary(did)
```

```
##
```

```
## Call:
## lm(formula = tot_traffic ~ treatment + after + treatment * after,
##     data = data1)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9924  -4874   2267   3928   6969
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    16992.2     1155.2  14.709 < 2e-16 ***
## treatment      139640.5     2310.5  60.438 < 2e-16 ***
## after          -677.2      2310.5  -0.293  0.771
## treatment:after -22721.4     4620.9  -4.917 1.27e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6003 on 44 degrees of freedom
## Multiple R-squared:  0.9904, Adjusted R-squared:  0.9897
## F-statistic: 1505 on 3 and 44 DF, p-value: < 2.2e-16
```

```
# Coefficients:
# Estimate Std. Error t value Pr(>|t|)
# (Intercept)    16992.2     1155.2  14.709 < 2e-16 ***
# treatment      139640.5     2310.5  60.438 < 2e-16 ***
# after          -677.2      2310.5  -0.293  0.771
# treatment:after -22721.4     4620.9  -4.917 1.27e-05 ***
# conclusion: treatment effect is significant
```

**Interpretation:** In the absence of running sponsored ads in google search engine, Bazaar.com on average would loose 22721 clicks per week.

This impact estimated over and above the control group behavior captures the true causality for the sponsored ads. This value is more accurate in comparison to the pre post estimate due to factoring of the control group behavior.

### 3. Given the effectiveness of branded keyword ads, what's the corrected RoI?

To estimate the Adjusted ROI, we need to estimate the proportion of users who were causally driven by the sponsored ads. In the earlier question, we arrived at the average traffic that were causally driven by ads. We also need determine the traffic that would have used organic search results in the absence of the ads to arrive at the proportion of traffic that were causally driven by ads.

```
# Traffic that would have landed in Bazaar.com even in absence of sponsored ads
did2 <- lm(avg_org ~ treatment + after + treatment * after, data=data1)
summary(did2)
```

```
##
## Call:
## lm(formula = avg_org ~ treatment + after + treatment * after,
##     data = data1)
##
## Residuals:
```

```
##      Min      1Q Median      3Q      Max
## -6244 -3212  2034   2807   3804
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    13043.2      706.8  18.454 < 2e-16 ***
## treatment     112255.9     1413.6  79.413 < 2e-16 ***
## after          -610.0     1413.6  -0.432  0.66818
## treatment:after  8545.0     2827.1   3.022  0.00417 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3673 on 44 degrees of freedom
## Multiple R-squared:  0.995, Adjusted R-squared:  0.9947
## F-statistic: 2914 on 3 and 44 DF, p-value: < 2.2e-16
```

**Interpretation:** In the absence of running sponsored ads in google, 8545 of the traffic would have used organic search results to land in Bazaar.com

Total clicks from sponsored ads (C) = clicks truly motivated by sponsored ads (A) + clicks by customers who would still visit Bazaar.com in the absence of sponsored ads (B)

A = 22721 (treatment effect estimated in did)

B = 8545 (treatment effect estimated in did2)

Proportion = A / (A+B) = 22721 / (22721+8545) = 0.7266999 = ~73%

The **Adjusted ROI** should be: (revenue per click \* Probability of click \* Proportion - cost per click) / cost per click

```
Adjusted_ROI = (21 * 0.12 * 0.7266999 - 0.6) / 0.6
Adjusted_ROI
```

```
## [1] 2.05214
```

Adjusted ROI = 205.214%

The corrected RoI from sponsored branded ads is determined to be ~205%. Though the adjusted ROI is lower than earlier estimate of 320%, it is still substantially high return for the amount invested. Hence, it is only appropriate to continue advertising through branded search ads.