

Project#1

Salem Almheiri

2/20/2022

R Markdown

```
library(tidyverse)
```

```
## — Attaching packages — tidyverse 1.3.1 —
```

```
## ✓ ggplot2 3.3.5      ✓ purrr 0.3.4
## ✓ tibble 3.1.6       ✓ dplyr 1.0.8
## ✓ tidyr 1.2.0        ✓ stringr 1.4.0
## ✓ readr 2.1.2        ✓ forcats 0.5.1
```

```
## — Conflicts — tidyverse_conflicts() —
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
library(jtools)
library(vtable)
```

```
## Loading required package: kableExtra
```

```
##
## Attaching package: 'kableExtra'
```

```
## The following object is masked from 'package:dplyr':
##
## group_rows
```

```
library(dplyr)
library(ggplot2)
library(wooldridge)
library(fixest)
library(haven)
library(lubridate)
```

```
##
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:base':  
##  
##    date, intersect, setdiff, union
```

```
library(purrr)
```

#The following code is used to remove the redundancy of having universities that share the same exact name

```
id_name_link <- read_csv("Lab3_Rawdata/id_name_link.csv")
```

```
## Rows: 3595 Columns: 3  
## — Column specification —————  
## Delimiter: ","  
## chr (1): schname  
## dbl (2): unitid, opeid  
##  
## i Use `spec()` to retrieve the full column specification for this data.  
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
id_name_link <- id_name_link %>%  
  group_by(schname) %>%  
  mutate(N = n()) %>%  
  filter(N == 1)
```

#Bind the multiple data sets on Google trends into one single data set

```
flist <- list.files('Lab3_Rawdata', pattern = 'trends', full.names = TRUE)  
trends <- flist %>%  
map(read_csv) %>%  
bind_rows()
```

```
## Warning: One or more parsing issues, see `problems()` for details
```

```
## Rows: 2825 Columns: 6  
## — Column specification —————  
## Delimiter: ","  
## chr (3): schname, keyword, monthorweek  
## dbl (3): schid, keynum, index  
##  
## i Use `spec()` to retrieve the full column specification for this data.  
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
## Warning: One or more parsing issues, see `problems()` for details
```

```
## Rows: 173949 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
## Rows: 168242 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
## Warning: One or more parsing issues, see `problems()` for details
```

```
## Rows: 139749 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
## Rows: 175729 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
## Warning: One or more parsing issues, see `problems()` for details
```

```
## Rows: 152884 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
## Warning: One or more parsing issues, see `problems()` for details
```

```
## Rows: 106965 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
## Rows: 115042 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
## Rows: 67236 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
## Rows: 306790 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
## Rows: 513 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
## Rows: 183816 Columns: 6
## — Column specification —————
## Delimiter: ","
## chr (3): schname, keyword, monthorweek
## dbl (3): schid, keynum, index
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
# Standardizing the Google Trends Data
```

```
trends <- trends %>%
  group_by(schname, keyword) %>%
  mutate(index_std = (index - mean(index, na.rm = TRUE))/sd(index, na.rm = TRUE))
```

```
#Upload the Scorecard data set properly
```

```
Scorecard <- read_csv("Lab3_Rawdata/Most+Recent+Cohorts+(Scorecard+Elements).csv")
```

```
## Rows: 7804 Columns: 122
```

```
## — Column specification —————
```

```
## Delimiter: ","
```

```
## chr (115): INSTNM, CITY, STABBR, INSTURL, NPCURL, LOCALE, HBCU, PBI, ANNHI, ...
```

```
## dbl (7): UNITID, OPEID, opeid6, HCM2, PREDDEG, CONTROL, CURROPER
```

```
##
```

```
## i Use `spec()` to retrieve the full column specification for this data.
```

```
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
#Change the name of the variable with the data on the median earnings of graduates ten y  
ears after graduation from "md_earn_wne_p10-REPORTED-EARNINGS" to "median.earnings.10yea  
rs"
```

```
names(Scorecard)[names(Scorecard) == "md_earn_wne_p10-REPORTED-EARNINGS"] <- "median.ear  
nings.10years"
```

```
#Change the name of the variable in the "Scorecard" data set from uppercase letters to l  
owercase ones to match the variables in the "id_name_link" data set
```

```
names(Scorecard)[names(Scorecard) == "UNITID"] <- "unitid"
```

```
names(Scorecard)[names(Scorecard) == "OPEID"] <- "opeid"
```

```
#The following code is used to match colleges as identified in the Scorecard data with c  
olleges as identified in the Google trends data that are then joined using the "id_name_  
link" data set
```

```
Joined <- id_name_link %>% left_join(Scorecard, by = "unitid")
```

```
Joined <- Joined %>% left_join(trends, by = "schname")
```

```
#Turn the "median.earnings.10years" variable from a string into a usable vector for comp  
utations
```

```
Joined <- Joined %>%
```

```
  mutate(median.earnings.10years = as.numeric(median.earnings.10years))
```

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

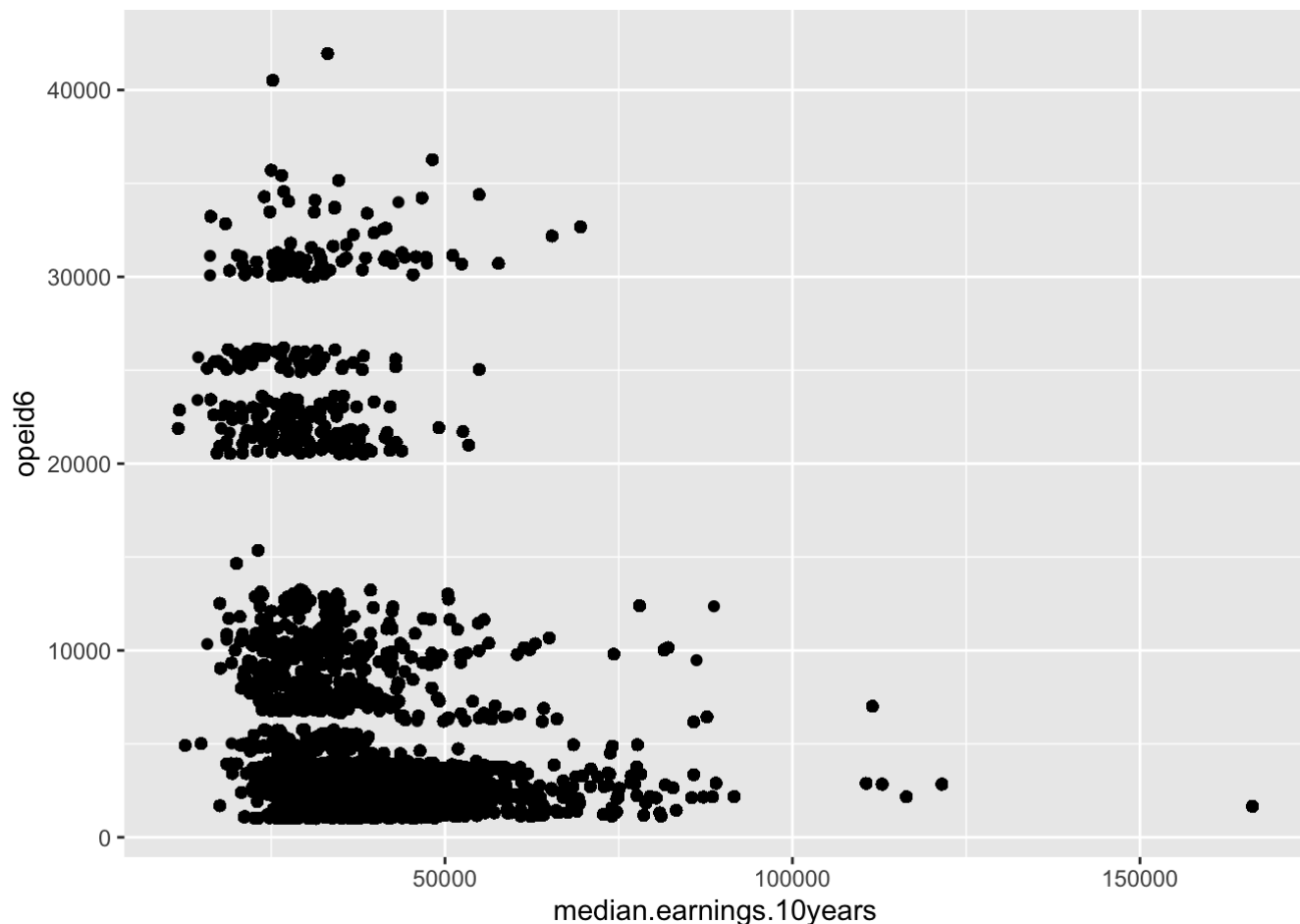
[illegible]


```
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
## Warning in mask$eval_all_mutate(quo): NAs introduced by coercion
```

#The following graph shows the median earning of a graduate 10 years after graduating on the x-axis and the variable "opeid6" on the y-axis (opeid6 is the best alternative to each college's name I found to visualize the data without having R crash; in this instance, the vertical axis, "opeid6", doesn't mean much, but it's a great way to visualize the data across the x-axis to see the data distribution). From this graph, we can clearly see that the majority of graduates earn less than \$100,000 ten years after graduating.

```
ggplot(Joined, aes(x = median.earnings.10years, y = opeid6)) + geom_point()
```

```
## Warning: Removed 94280 rows containing missing values (geom_point).
```



#However, this is not the full story. Our study will focus on colleges that predominantly grant bachelor's degrees, so we need to see how the data looks like when we look at only such colleges. I will create a sub-sample first using the following code:

```
Subsample <- Joined %>% filter(PREDDEG == "3")
```

#Find the median of the earnings of graduates 10 years after graduating in the sub-sample, which only includes colleges that predominantly grant bachelor's degrees

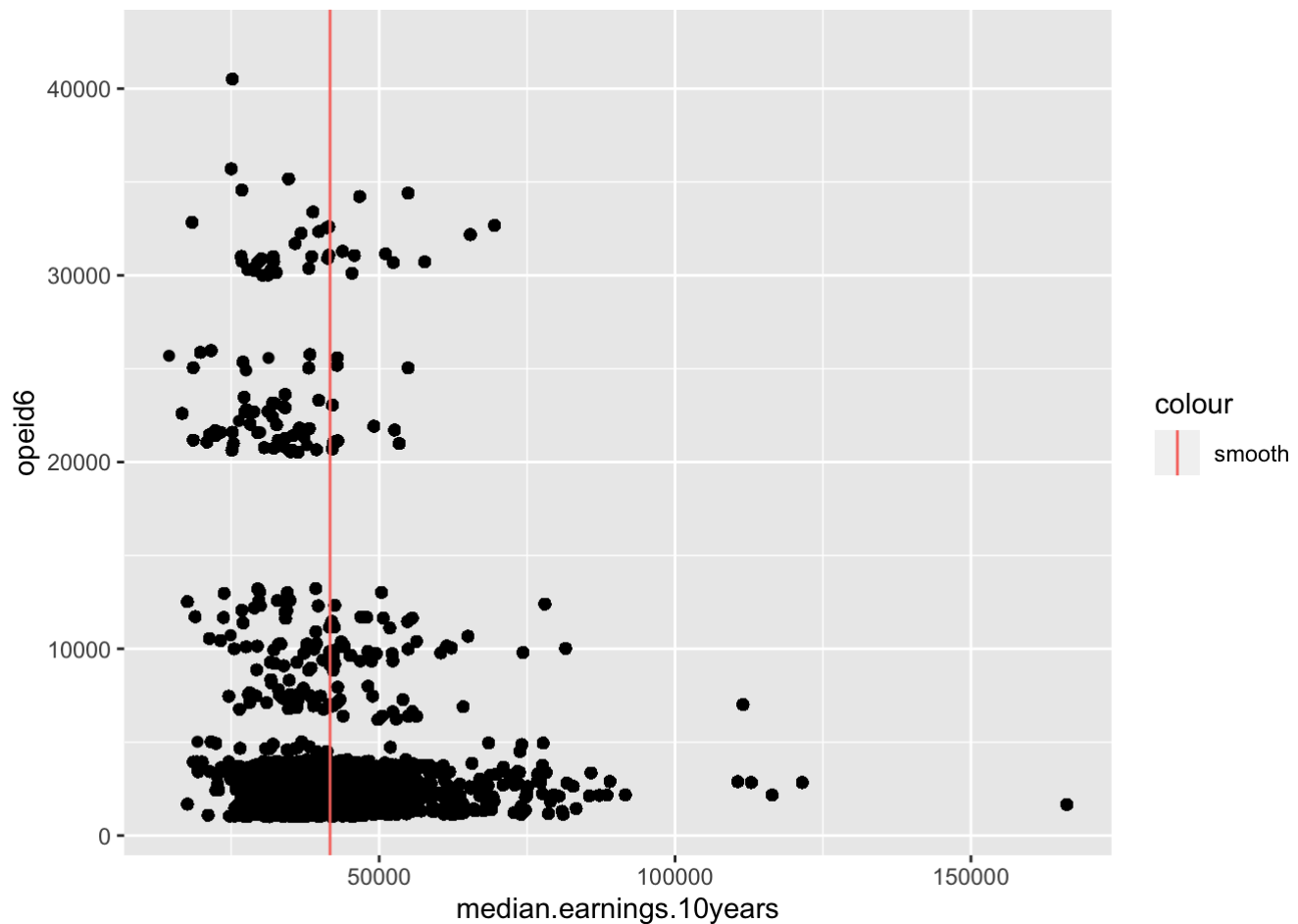
```
Subsample %>% pull(median.earnings.10years) %>% median(na.rm = TRUE)
```

```
## [1] 41700
```

#This value will be the cut off to differentiate between high-earning and low-earning colleges (among those that predominantly grant bachelor's degrees). This means if the college's graduates are earning \$41,700 or above ten years after graduating college, they graduated from a high-earning college. Similarly, if a college's graduates are earning below \$41,700 ten years after graduating college, they graduated from a low-earning college. This can be better seen in the following graph:

```
ggplot(Subsample, aes(x = median.earnings.10years, y = opeid6)) + geom_point() + geom_vline(aes(xintercept = 41700, colour = "smooth"))
```

```
## Warning: Removed 67887 rows containing missing values (geom_point).
```



#Here, I will create dummy variables that I plan to use in my regression

#The following adds the dummy variable "after_scorecard" in the sub-sample that I created earlier. When this variable is true, it refers to the period after early September 2015, which is when the scorecard data was released.

```
Subsample <- Subsample %>%
  mutate(date = str_sub(monthorweek, 1, 10)) %>%
  mutate(date = ymd(date)) %>%
  mutate(after_scorecard = date >= ymd('2015-09-01'))
```

#The following adds the dummy variable "high_earning" in the sub-sample. When this variable is true, it refers to high-earning colleges (that predominantly grant bachelor's degrees), where graduates are earning \$41,700 or above ten years after graduating college

```
Subsample <- Subsample %>%
  mutate(high_earning = median.earnings.10years >= 41700)
```

#The following is my regression equation. I chose this regression design because it will allow me to see how the before/after change in student interest was different for high-earning colleges as compared to low-earning ones

```
Regression <- lm(index_std ~ after_scorecard*high_earning, data = Subsample)
```

#The following code will show the OLS coefficient terms in a table

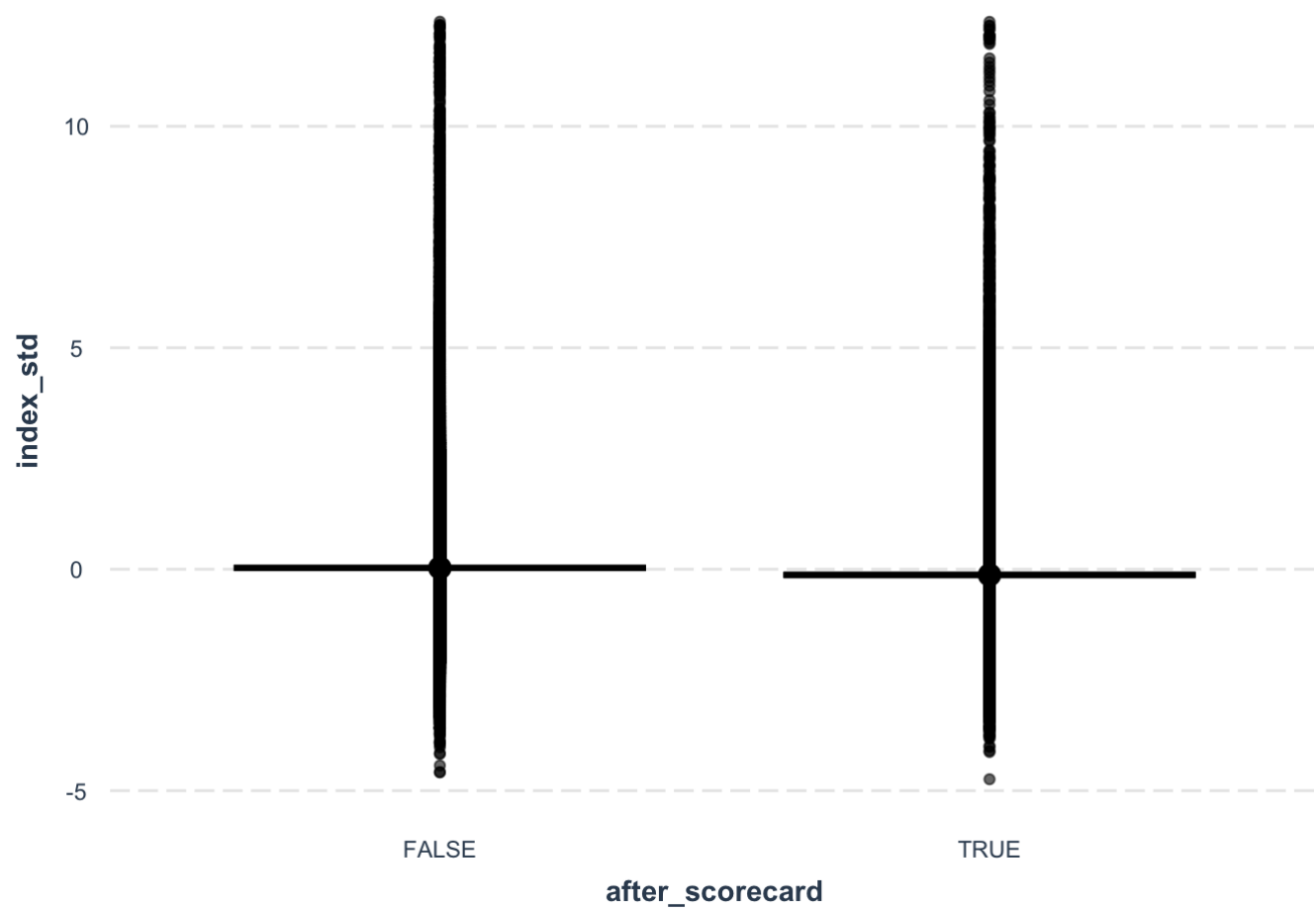
```
export_summs(Regression)
```

	Model 1
(Intercept)	0.03 *** (0.00)
after_scorecardTRUE	-0.16 *** (0.00)
high_earningTRUE	0.01 *** (0.00)
after_scorecardTRUE:high_earningTRUE	-0.06 *** (0.01)
N	910198
R2	0.01

*** p < 0.001; ** p < 0.01; * p < 0.05.

#The following code will plot the regression on a graph

```
effect_plot(Regression,
  pred = after_scorecard,
  plot.points = TRUE)
```



#Interpretations of the coefficient terms table

#Among low-earning colleges that predominantly grant bachelor's degrees, before the release of the scorecard, student interest (as measured by Google searches for keywords associated with those colleges) was 0.03 units higher.

#Among low-earning colleges that predominantly grant bachelor's degrees, the release of the scorecard decreased student interest (as measured by Google searches for keywords associated with those colleges) by 0.16 units.

#Among high-earning colleges that predominantly grant bachelor's degrees, before the release of the Scorecard, student interest (as measured by Google searches for keywords associated with those colleges) was 0.01 units higher.

#Among high-earning colleges that predominantly grant bachelor's degrees, the release of the Scorecard decreased student interest (as measured by Google searches for keywords associated with those colleges) by 0.05 units.

#So, it seems like the introduction of the Google Scorecard decreased student interest (as measured by Google searches for keywords associated with those colleges) overall. Nonetheless, the magnitude of decrease in student interest we notice in low-earning colleges is comparatively much higher than it is in high-earning college. As such, we can conclude that the release of the college scorecard at the start of September 2015 did shift student interest to high-earning colleges relative to low-earning ones (as proxied by Google searches for keywords associated with those colleges) among colleges that predominantly grant bachelor's degrees.