

# Data Exploration: Symbolic Politics

Luke Kolar

October 21, 2021

In this Data Exploration assignment we will explore Reny and Newman's (2021) finding that opinions towards the police and about the level of discrimination faced by Black Americans were impacted by the spread of protests in the wake of the killing of George Floyd. You will recreate, present, and assess those claims as well as creating your own regression models to test which attitudes change and when.

If you have a question about any part of this assignment, please ask! Note that the actionable part of each question is **bolded**.

## Opinion Mobilization: The George Floyd Protests

### Data Details:

- File Name: `RN_2001_data.RData`
- Source: These data are from Reny and Newman (2021).

Variable Name	Variable Description
<code>race_ethnicity</code>	Race or ethnicity. Levels labelled in data: 1-White, 2-Black or AfAm, 3-American Indian or Alaskan Native, 4 through 14- Asian or Pacific Islander (details in labels), and 15-Some other race
<code>hispanic</code>	Of Hispanic, Latino, or Spanish origin. Levels labelled in data: 1-Not Hispanic, 2-15 Hispanic of various origins
<code>day_running</code>	Day relative to onset of George Floyd protests (day 0)
<code>age</code>	Respondent's age
<code>female</code>	Binary indicator variable: 1 if respondent female, 0 otherwise
<code>college</code>	Binary indicator variable: 1 if respondent attended college, 0 otherwise
<code>household_income</code>	Household pre-tax income ranging from 1 (less than \$15,000) to 24 (more than \$250,000). Details for other levels in labels.
<code>pid7</code>	Party identification on a seven point scale with strong, weak, lean: 1-Strong Democrat to 7-Strong Republican with 4-Independent.
<code>ideo5</code>	Ideological self placement: 1-Very liberal, 2-Liberal, 3-Moderate, 4-Conservative, 5-Very Conservative
<code>vote_clinton</code>	Indicator variable for whether the respondent said they voted for Clinton in the 2016 presidential election
<code>group_favorability_the_police</code>	Favorability towards the police: 1-Very favorable, 2-Somewhat favorable, 3-Somewhat unfavorable, 4-Very unfavorable
<code>discrimination_blacks</code>	Perceptions of the level of discrimination in US faced by Blacks: 1-None at all, 2-A little, 3-A moderate amount, 4-A lot, 5-A great deal

Variable Name	Variable Description
day	The date the respondent took the survey
group_fav_white_black	The difference in respondents favorability towards Blacks subtracted from their favorability towards whites (each on four point scale). Ranges from -3 to 3.
racial_attitudes_generations	Agreement with the statement that generations of slavery and discrimination have made it difficult for Blacks to work their way out of the lower class: 1-Strongly Agree to 5-Strongly Disagree
interest	Degree to which respondent claims to follow politics: 1-Most of the time, 2-Some of the time, 3-Only now and then, 4-Hardly at all
group_favorability_jews	Favorability towards Jews: 1-Very favorable, 2-Somewhat favorable, 3-Somewhat unfavorable, 4-Very unfavorable
group_favorability_whites	Favorability towards whites: 1-Very favorable, 2-Somewhat favorable, 3-Somewhat unfavorable, 4-Very unfavorable
group_favorability_evangelicals	Favorability towards evangelicals: 1-Very favorable, 2-Somewhat favorable, 3-Somewhat unfavorable, 4-Very unfavorable
group_favorability_socialists	Favorability towards socialists: 1-Very favorable, 2-Somewhat favorable, 3-Somewhat unfavorable, 4-Very unfavorable
protest	Indicator variable if survey respondent lived in area that would at any point have a BLM protest in the wake of the killing of George Floyd
n_protests	Number of eventual BLM protests in area where resident lived

```
# load the data containing the tibble protest_df
load('RN_2001_data.RData')
```

```
#Note that the data is saved in the form of a tibble, a special table using the dplyr package that has
head(protest_df$race_ethnicity)
```

```
## <labelled<double>[6]>: What is your race? Provided by LUCID.
## [1] 6 1 1 2 1 1
##
## Labels:
##   value          label
##     1             White
##     2   Black, or African American
##     3 American Indian or Alaska Native
##     4       Asian (Asian Indian)
##     5       Asian (Chinese)
##     6       Asian (Filipino)
##     7       Asian (Japanese)
##     8       Asian (Korean)
##     9       Asian (Vietnamese)
##    10       Asian (Other)
##    11 Pacific Islander (Native Hawaiian)
##    12   Pacific Islander (Guamanian)
##    13   Pacific Islander (Samoan)
##    14   Pacific Islander (Other)
##    15       Some other race
##   777   Not asked in this wave
```

```
head(protest_df$household_income)
```

```
## <labelled<double>[6]>: What is your current annual household income before taxes? Provided by L...
## [1] 21  8  7  1 NA  1
##
## Labels:
##   value          label
##    1      Less than $14,999
##    2    $15,000 to $19,999
##    3    $20,000 to $24,999
##    4    $25,000 to $29,999
##    5    $30,000 to $34,999
##    6    $35,000 to $39,999
##    7    $40,000 to $44,999
##    8    $45,000 to $49,999
##    9    $50,000 to $54,999
##   10    $55,000 to $59,999
##   11    $60,000 to $64,999
##   12    $65,000 to $69,999
##   13    $70,000 to $74,999
##   14    $75,000 to $79,999
##   15    $80,000 to $84,999
##   16    $85,000 to $89,999
##   17    $90,000 to $94,999
##   18    $95,000 to $99,999
##   19  $100,000 to $124,999
##   20  $125,000 to $149,999
##   21  $150,000 to $174,999
##   22  $175,000 to $199,999
##   23  $200,000 to $249,999
##   24    $250,000 and above
##  777 Not asked in this wave
```

## Question 1

As usual it is important to first examine the structure of the data. What are the two main outcome variables of interest to Reny and Newman? How were they measured and how are they coded in the data? What was the treatment? **Take a look at the data and determine which are the two outcome variables of interest. Observe the scale of each.**

```
glimpse(protest_df)
```

```
## Rows: 378,507
## Columns: 22
## $ race_ethnicity      <dbl+lbl> 6, 1, 1, 2, 1, 1, 1, 1, 1, 1, 1...
## $ hispanic            <dbl+lbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
## $ day_running         <dbl> -315, -315, -315, -315, -315, -315,...
## $ age                 <dbl> 37, 45, 24, 26, 60, 55, 46, 37, 60,...
## $ female              <dbl> 1, 0, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1,...
## $ college             <dbl> 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1,...
## $ household_income    <dbl+lbl> 21,  8,  7,  1, NA,  1, 19, 20,...
## $ pid7               <dbl+lbl>  2,  1,  1,  3,  7,  2,  7, NA,...
```

```
## $ ideo5 <dbl+lbl> 4, 2, 2, 3, 5, 2, 5, 3, 5, 1, 4...
## $ vote_clinton <dbl> 0, 1, 1, 1, 0, 1, 0, 0, 0, 1, 0, 0,...
## $ group_favorability_the_police <dbl+lbl> 2, 2, 2, 4, 1, 3, 1, 2, 2, 3, 1...
## $ discrimination_blacks <dbl> 3, 1, 5, 5, 1, 5, 2, 4, 3, 5, 3, 3,...
## $ day <date> 2019-07-18, 2019-07-18, 2019-07-18...
## $ group_fav_white_black <dbl> 0, 3, -1, -2, 2, 0, 0, -2, 0, -2, 1...
## $ racial_attitudes_generations <dbl+lbl> 4, 5, 2, 1, 5, 2, 5, 3, 2, 1, 5...
## $ interest <dbl+lbl> 2, 1, 2, 2, 1, 2, 1, 4, 3, 2, 1...
## $ group_favorability_jews <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA,...
## $ group_favorability_whites <dbl> 3, 4, 3, 2, 4, 2, 3, 1, 2, 2, 4, 3,...
## $ group_favorability_evangelicals <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA,...
## $ group_favorability_socialists <dbl> 2, 1, 3, NA, 4, 3, 4, 3, NA, 1, 4, ...
## $ protest <dbl> 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1,...
## $ n_protests <dbl> 7, 0, 2, 4, 2, 4, 0, 7, 1, 7, 3, 6,...
```

**ANSWER:** *It seems that the two main outcome variables for Reny and Newman (2021) are group\_favorability\_the\_police and day\_running. The former tells us how respondents feel about the police, and the latter gives us a relative timeframe through which to investigate this - before and after the protests.*

##Question 2

###Part a R has a special 'date' class for storing and manipulating dates as seen below. Date variables can conveniently be logically compared and arithmetically manipulated. Using the day variable find out how many days the dataset spans. **First check using the code below that the day variable is of the class 'date'.** Next subtract the latest day in the sample from the first day to calculate the timespan covered by the dataset. Hint: functions like max() and min() work for date variables too!

```
class(protest_df$day)
```

```
## [1] "Date"
```

```
days <- pull((protest_df %>% summarize(max(day))) - (protest_df %>% summarize(min(day))))
days
```

```
## Time difference of 419 days
```

###Part b On what date is the treatment said to have occurred? **Find the date for which the day\_running variable is 0. Then modify the code below to add a variable to each row for whether or not the observation was before or after treatment.**

*#Change the object to be the date of the protest spread, remember to put it in quotes if you copy/paste*

```
date <- protest_df %>%
  filter(day_running == 0) %>%
  head(1) %>%
  pull(day)

protest_df_bydate <- protest_df %>%
  mutate(before = ifelse(day < as.Date(date), 1, 0))
```

### Question 3

###Part a Compare the average for each outcome variable before and after the onset of the protests. Are the differences statistically significant? **Calculate the outcome variable means for before and after treatment. Conduct a test as to whether the differences in means are statistically significant. Hint: you can use either the `t.test()` function or `difference_in_means()` from the `estimatr` package**

```
t.test(protest_df$group_favorability_the_police, protest_df$before)
```

```
##
## One Sample t-test
##
## data: protest_df$group_favorability_the_police
## t = 1177.7, df = 344699, p-value < 2.2e-16
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
##  2.038667 2.045464
## sample estimates:
## mean of x
##  2.042066
```

**ANSWER:** *Given the 95% confidence interval does not include zero, there is statistically significant evidence that the true difference in means between favorability towards police before and after the protests is greater than 0.*

###Part b It might be that the period before and after the treatment was different in ways in addition to the onset of the protests. Use the same procedure as above to check for differences between two means of a survey response measuring favorability towards a group besides the police. **Calculate the means from before and after the treatment and conduct a test of statistical significance of the difference for another measure of group favorability that was recorded in the survey (e.g. evangelicals, Jews, socialists, or whites).** Is there also a substantive or statistically significant difference on that variable? Should that change our confidence in attributing the opinion changes found in part a to the George Floyd protests?

### Question 4

###Part a In order to create figures similar to the panels in Figure 2 in Reny and Newman (2021) we must first manipulate the data to be more usable. If we intend to graph the average of each outcome variable for each day, on what variable should we group the data using `group_by`? **Create a new object that is the data split out by the appropriate group and producing the average for each of the two outcome variables for each day. Also be sure to preserve an indicator for whether the observations are from before or after the spread of the protests.**

###Part b Graph the results for the entire sample. **Graph the results for the entire sample for both outcome variables by day. Include a vertical line demarcating when the protests started to spread. Does there appear to be a shift in the outcome variables from before to after the protests began to spread?**

###Part c It might be useful to more clearly illustrate the differences in the trend lines before and after the protests began. **Modify the code below to include a separate line of best fit for before and after the protests began. Does the trend line align with your previous reading of the graph? Remember to add a vertical line demarcating for the onset of treatment.**

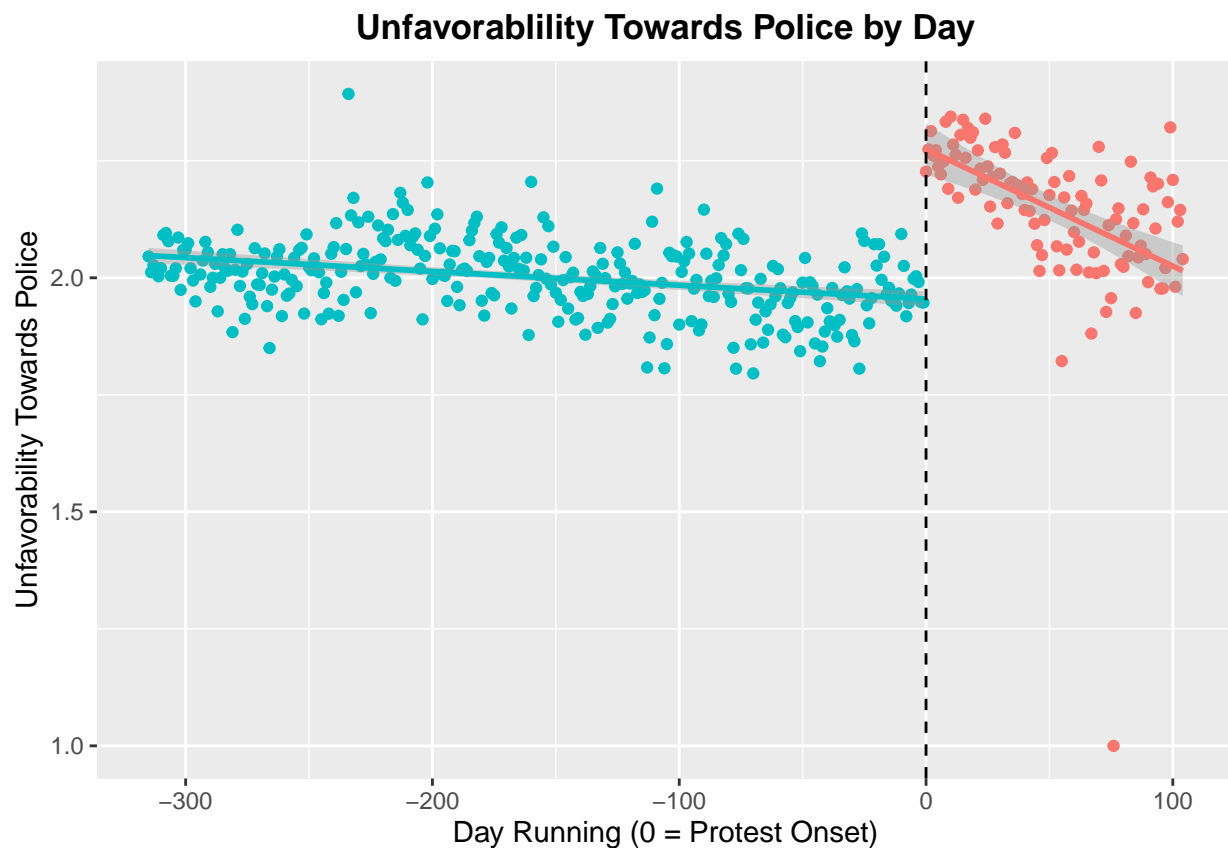
```

protest_df %>%
  group_by(day_running) %>%
  summarize(avg_day = mean(group_favorability_the_police, na.rm = TRUE)) %>%
  mutate(before = ifelse(day_running >= 0, FALSE, TRUE)) %>%
  ggplot(aes(x = day_running, y = avg_day, color = before)) + geom_point() +
  geom_smooth(aes(color = before), method = "lm") +
  labs(title = "Unfavorablility Towards Police by Day", x = "Day Running (0 = Protest Onset)",
       y = "Unfavorability Towards Police") +
  geom_vline(aes(xintercept = 0), linetype = "dashed") +
  theme(legend.position = "none",
        plot.title = element_text(hjust = 0.5, face = "bold"))

```

## 'summarise()' ungrouping output (override with '.groups' argument)

## 'geom\_smooth()' using formula 'y ~ x'



## Question 5

### Part a The attitudes in question are no doubt highly influenced by the respondent's race and ethnicity. How do the graphs from question 4 differ for white and Black respondents. **Subset the data to include only white respondents and recreate the graphs from part c of question 4. Do the same with the data from only Black respondents. How do these differ from each other? Hint: Be careful when subsetting white responses to not also include Hispanic responses.**

###Part b As we have learned partisanship heavily influences how people take in and process new information. **Split the sample into Democrats, Republicans and independents and use them to produce the same graphs as part a (either all in the same figure or separate). Compare both the level and the trends for each party affiliation. What could this imply about how partisanship affects processing?**

##Question 6:

###Part a The graphs in questions 4 and 5 indicate that the effects dissipate as time progresses past the onset of the protests. **Explain why that might be the case? What does this indicate about whether or not attitudes towards the police are symbolic or not?**

###Part b One way to look at the effect decay is to bin the post-protest data and compare averages. **Split the post-protest data into however many groups you choose and compare the period directly after the protest with the latest period in the data. What are the differences in means for the outcomes?**

## Question 7

###Part a What are some reasons we might be unconvinced by the comparison of aggregate survey results from a time before and after an event? Do you think they apply here?

###Part b There is often a problem in conducting surveys of non-response bias. That is, the people who answer surveys may differ from the people who do not answer surveys and the differences may vary over time. This is especially damaging to inference when non-response is correlated with the outcomes being measured. For example after a series of damaging headlines supporters of a politician may be less willing to answer phone surveys about that politician. As a result we would potentially observe an exaggeration of the negative effects of the scandal on a politician's polled approval rating. Test whether this is the case in the Reny and Newman data. **Test whether there is balance between the respondents before and after the onset of the protests along two demographic traits that you would expect to correlate with the measured responses to the outcome variables.**

###Part c Racial resentment is often considered a symbolic attitude in strength and consistency. Examine the before and after levels of racial resentment as measured by the question from the racial resentment scale about the impact of generations of slavery and discrimination (racial\_attitudes\_generations). **Graph the average racial\_attitudes\_generations (remember the direction of how it is coded!) by day like other outcome variables. Does it behave like the other outcome variables? Does the data support that racial attitudes are symbolic attitudes?**

## Question 8: Data Science Question

###Part a Run an initial regression examining the relationship between favorability towards the police, party, and treatment. **Run a regression examining party and the onset of the protests' effect on favorability towards the police. Interpret the results**

```
protest_df_edit <- protest_df %>%
  mutate(treatment_after = ifelse(day_running >= 0, 1, 0))

mod = lm(group_favorability_the_police ~ treatment_after + pid7, data = protest_df_edit)

stargazer(mod, type = "text")
```

```
##
## =====
```

```
##                               Dependent variable:
##                               -----
##                               group_favorability_the_police
##                               -----
## treatment_after              0.169***
##                               (0.004)
##
## pid7                         -0.127***
##                               (0.001)
##
## Constant                     2.478***
##                               (0.003)
##
## -----
## Observations                 326,797
## R2                           0.088
## Adjusted R2                  0.088
## Residual Std. Error          0.966 (df = 326794)
## F Statistic                   15,807.950*** (df = 2; 326794)
## =====
## Note:                        *p<0.1; **p<0.05; ***p<0.01
```

**ANSWER:** The regression suggests that there is a statistically significant relationship at the  $p < 0.01$  level between unfavorability towards the police and both party identification and party identification (scaled 1 to 7). The coefficient for the treatment, or whether or not the protests occurred, is 0.169, which means that with all other variables held constant, on average, we can expect views after the protests to be higher by 0.169. The coefficient for the party identification variable is -0.127, which means that with all other variables held constant, a one-point increase on the pid7 scale (which indicates a move toward Republicanism) results in an expected -0.127 decrease in police unfavorability.

###Part b The above functional form probably does not accurately model the relationship of all the relevant covariates in the dataset. What functional form would you recommend using and why? What covariates would you add? Is there need for an interaction term? **Run a regression of your specification and interpret the results. Justify your choices in modeling.**

```
protest_df_edit2 <- protest_df_edit %>%
  mutate(is_black = ifelse(race_ethnicity == 2, 1, 0))

mod = lm(group_favorability_the_police ~ treatment_after + ideo5*is_black, data = protest_df_edit2)

stargazer(mod, type = "text")
```

```
##
## =====
##                               Dependent variable:
##                               -----
##                               group_favorability_the_police
##                               -----
## treatment_after              0.158***
##                               (0.004)
##
## ideo5                       -0.248***
##                               (0.002)
##
```



```
## is_black          0.213***
##                  (0.015)
##
## ideo5:is_black     0.150***
##                  (0.005)
##
## Constant          2.676***
##                  (0.005)
##
## -----
## Observations      344,101
## R2                 0.119
## Adjusted R2       0.119
## Residual Std. Error 0.955 (df = 344096)
## F Statistic       11,650.980*** (df = 4; 344096)
## =====
## Note:              *p<0.1; **p<0.05; ***p<0.01
```

**ANSWER:** I chose to replace the `pid7` variable with one that measures ideology: the `ideo5` variable. I wanted to keep my model simple, but strong. So, I also added a variable, `is_black`, which is coded 1 if the respondent is Black and 0 otherwise. I included an interaction variable between the two as well to account for conservation Blacks who may be pro-police. I toyed with many variables, including household income - the coefficient for this was too small to justify including, as was the case with college education level, proximity to the protests, and interest in politics. Variables like attitudes toward Blacks or Whites were highly correlated with unfavorability toward police, but I believe these may be correlated with some of the more fundamental forces represented in my model. Ultimately, my regression model was strong, with all coefficients displaying a statistically significant relationship at the  $p < 0.01$  level with unfavorability towards police. The treatment had a coefficient of 0.158, which means we can expect a 0.158 increase in unfavorability if the protests occurred (all other variables held constant, and at the 99% confidence level). For a one-point increase on the ideology scale, a step toward conservatism, we can expect a -0.248 change in unfavorability toward the police. On average, we can expect that Blacks have a 0.213 higher unfavorability “score” when compared to non-Blacks. And if a respondent is Black, the interaction term says that for every one-point increase on the ideology scale, we can expect a 0.150 increase in unfavorability toward the police. This last coefficient is the most interesting, as it suggests that more conservative Blacks have less favorable views toward the police.

### Part c Linear models are not well suited for bounded ordinal responses. Instead ordinal logit or probit models are frequently employed in order to capture a) that the outcomes are restricted to a scale (in the case of police unfavorability 1-4) and b) that the differences between different rungs on the scale are not necessarily equivalent (going from very unfavorable to somewhat unfavorable is not necessarily the same difference as going from somewhat unfavorable to somewhat favorable). **Using the code below from the MASS package run an ordinal probit model using the same model as part b. How do the coefficients differ from part b?**

```
library(MASS)
```

```
##
## Attaching package: 'MASS'

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

```

select <- dplyr::select

protest_df_edit3 <- protest_df_edit2 %>%
  mutate(group_favorability_the_police = as.factor(group_favorability_the_police),
         treatment_after = as.factor(treatment_after),
         ideo5 = as.factor(ideo5),
         is_black = as.factor(is_black))

polr(data = protest_df_edit3, formula = group_favorability_the_police ~ treatment_after + ideo5*is_black,
     method = "probit")

## Call:
## polr(formula = group_favorability_the_police ~ treatment_after +
##       ideo5 * is_black, data = protest_df_edit3, method = "probit")
##
## Coefficients:
## treatment_after1      ideo52      ideo53      ideo54
##      0.1690022    -0.2372634    -0.4688086    -0.8843135
##      ideo55      is_black1 ideo52:is_black1 ideo53:is_black1
##     -1.1527623      0.3292775      0.1920383      0.3833238
## ideo54:is_black1 ideo55:is_black1
##      0.5704331      0.6934550
##
## Intercepts:
##      1|2      2|3      3|4
## -0.7511156  0.1747168  0.8495425
##
## Residual Deviance: 846165.48
## AIC: 846191.48
## (34406 observations deleted due to missingness)

```

**ANSWER:** Most coefficients are significantly larger than those in part b. This makes sense - if we are able to analyze the impact of, say, ideology at different levels of ideology, rather than a one-point increase from any level, we're able to better pinpoint the relationship between that specific level and the dependent variable. In this case, there's an especially stronger negative relationship for each step higher on the ideology scale toward more conservative views. Interesting, the interaction variable shows the opposite effect for higher levels of ideology if the respondent is Black - coefficients are larger in the positive direction.