

# w203 Lab 1: Comparing Means

Srishti Mehra, David Djambazov, and Andi Morey Peterson

10/17/2020

## The Data

The American National Election Studies (ANES) conducts surveys of voters in the United States. While its flagship survey occurs every four years at the time of each presidential election, ANES also conducts pilot studies midway between these elections. We will be using this data to ask five (5) questions about the respondents:

1. Do US voters have more respect for the police or for journalists?
2. Are Republican voters older or younger than Democratic voters?
3. Do a majority of independent voters believe that the federal investigations of Russian election interference are baseless?
4. Was anger or fear more effective at driving increases in voter turnout from 2016 to 2018?
5. (Student Choice) Do Hillary Clinton voters and Trump voters view the US income gap differently?

## General Study Comments (applicable for all questions)

Since all questions draw from the same study sample, it is useful to state comments here that we can refer to for all questions. For almost any test we run, we will need to determine if the data is i.i.d. and if the respondents to the survey accurately represent the average U.S. voter that we can generalize accordingly.

*Independence* - The ANES study was conducted on the Internet using the YouGov panel. The YouGov panel consists of a large and diverse set of over a million volunteer respondents. Respondents were selected from the YouGov panel by sample matching. Given this methodology description from the ANES survey, we conclude that it is extremely unlikely for large clusters of people who somehow know each other to be present in the dataset, so we can safely assume that each respondent is independent from any other.

*Identically Distributed* - Once a person has taken the survey, they cannot take the survey again; so the distribution for the next “draw” is changed. But the change in the population distribution for the next draw is so small, we can safely ignore this effect.

*Generalizability* - Because this is a modern, paid, opt-in survey, the sample data will only include individuals who have the propensity or the financial motivation to complete online surveys. However, the financial impact is small, 21-50 cents for this 30 minute survey (see the ANES User Guide Code Book). In addition, the survey provided weights which the authors recommend using when making inferences about the general population of U.S. adult citizens.

Given these, we can assume the iid assumption is valid and use the weights to help us on questions in which we are concerned about generalizing to the population. (One concern/gap worth mentioning – the dataset

does not account for people who are ineligible to vote due to a felony. Such respondents wouldn't have voted because they are unable to do so, not for any other reason).

*Confidence Interval* - All tests will use a 95% confidence as standard practice.

*Voting Population* - Nearly all the questions ask about voters. The survey provides a few sample questions about the respondent to determine if that respondent was actually a voter. Variables *turnout18* and *turnout18ns* can be used to determine if they were a voter. For our analysis, we will only consider responses that have value 1, 2, or 3 for the variable *turnout18*.

Let's look at how many observations are in those categories:

```
paste("Number of NA in turnout18: ", length(which(is.na(A$turnout18))))
```

```
## [1] "Number of NA in turnout18: 0"
```

```
paste("Number of NA in turnout18ns: ", length(which(is.na(A$turnout18ns))))
```

```
## [1] "Number of NA in turnout18ns: 0"
```

```
paste("Definitely Voted: ", sum(A$turnout18 <= 3))
```

```
## [1] "Definitely Voted: 1842"
```

```
paste("Not completely sure or Probably did vote: ", sum(A$turnout18 == 5 & A$turnout18ns == 1))
```

```
## [1] "Not completely sure or Probably did vote: 18"
```

```
paste("Number of Duplicate Voters for the main dataset: ", length(which(duplicated(A))))
```

```
## [1] "Number of Duplicate Voters for the main dataset: 0"
```

```
A$voted2018 <- ifelse((A$turnout18<=3), 1, 0)
```

About 1% are not sure AND “probably” voted. We argue that because this is a fairly small number and being uncertain about having voted in an election that has just taken place can be reasonably viewed as grounds for exclusion from the population of US 2018 voters. We will use this population for all questions in this lab asking specifically about “voters”.

We also did a check to confirm there are no duplicate rows in the set of *voted18*.

# Research Questions

## Question 1: Do US voters have more respect for the police or for journalists?

### Introducing the topic

Concept: To understand if US voters have different amount of respect police and journalists.

Operationalization: From a list of data points collected in the ANES 2018 Pilot Study, two are respondent's ratings of the police and of journalists in variables `ftpolicy` and `ftjournal` respectively. The questions to collect these are "How would you rate the police?" and "How would you rate journalists?". These are the closest variables to our question available since they describe how the sample population feels about the police and journalists respectively. The answers to these questions are collected on a rating scale between 0 to 100.

To get the subset of voters, we will use the `voted2018` variable described in the introduction.

Concerns or Gaps:

1. The data collected is a sentiment rating of police and journalists in general, not specifically addressing "respect". We have to be cautious while interpreting our conclusions since we are inferring about respect from a total rating.

### Exploratory data analysis (EDA) of the relevant variables

Variables: `ftpolicy`, `ftjournal`

Types: Both `ftpolicy` and `ftjournal` are ordinal variables.

Number of Entries

```
library(tidyverse)
usVoters <- A %>% filter(A$voted2018 == 1)
paste("Number of Entries:", length(usVoters$caseid))
```

```
## [1] "Number of Entries: 1842"
```

There are no Null Values for `ftpolicy` and `ftjournal`

```
paste("Number of Null Values for ftpolicy:", sum(is.na(usVoters$ftpolicy)))
```

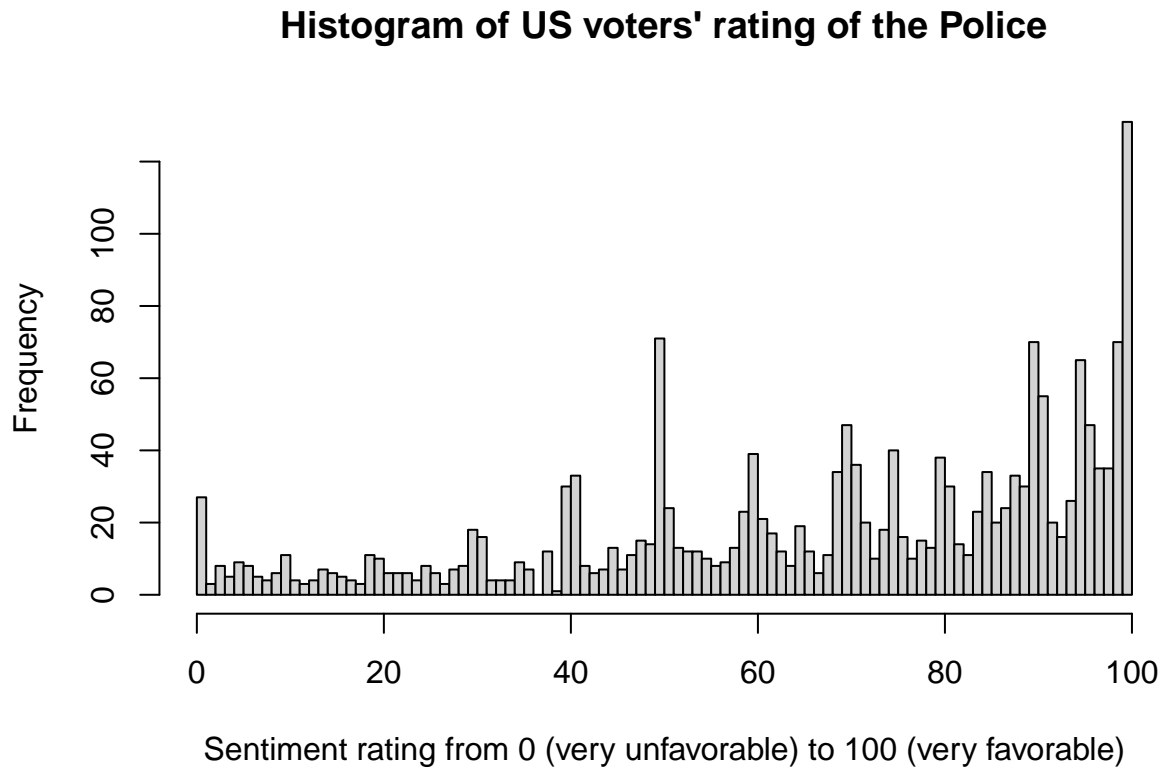
```
## [1] "Number of Null Values for ftpolicy: 0"
```

```
paste("Number of Null Values for ftjournal:", sum(is.na(usVoters$ftjournal)))
```

```
## [1] "Number of Null Values for ftjournal: 0"
```

We will draw a histogram of both `ftpolicy` and `ftjournal` to observe the distribution and any outliers.

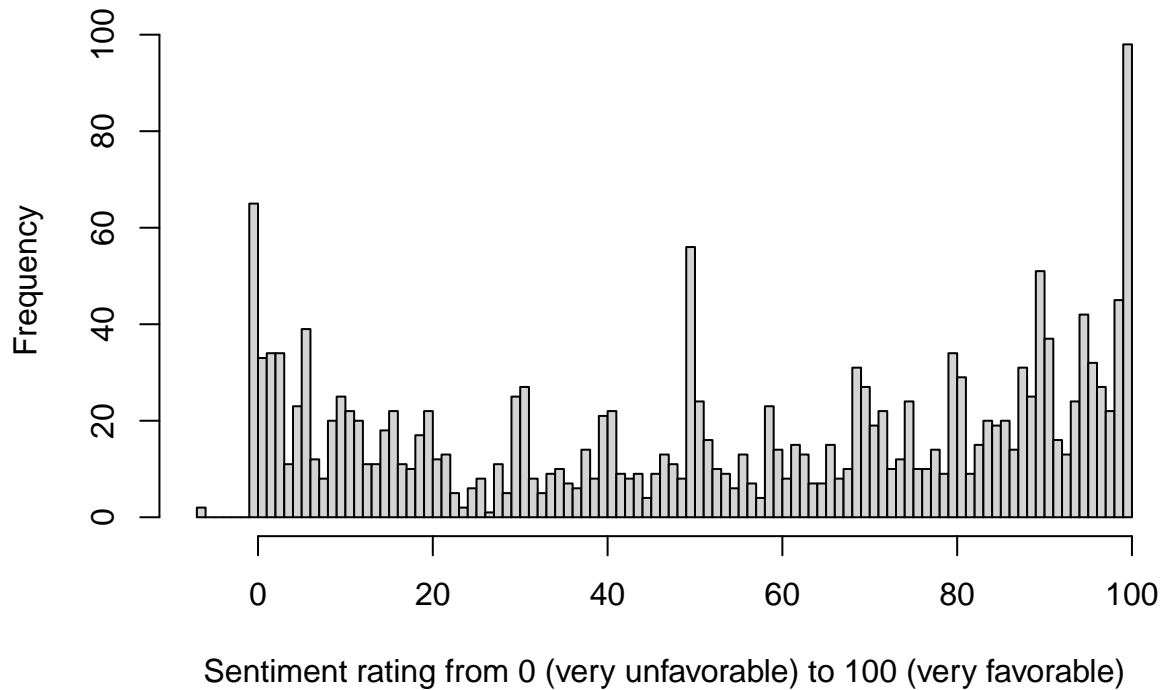
```
hist(usVoters$ftpolice, breaks = 100,
     main = "Histogram of US voters' rating of the Police",
     xlab = "Sentiment rating from 0 (very unfavorable) to 100 (very favorable)")
```



> The distribution for `ftpolice` seems a little skewed towards the higher end, with upticks in the center.

```
hist(usVoters$ftjournal, breaks = 100,
     main = "Histogram of US voters' rating of the Journalists",
     xlab = "Sentiment rating from 0 (very unfavorable) to 100 (very favorable)")
```

## Histogram of US voters' rating of the Journalists



The distribution for `ftjournal` is less skewed than that of `ftpolic`.

### Based on our EDA, selecting an appropriate hypothesis test

1. Since we are comparing two samples, it will be a two-sample test.
2. Since the variables are ordinal, we will use a non-parametric test.
3. Since the two variables we are looking at are answered by the same person, we can do a paired test. For the above reasons, we narrow our choice of test down to Sign Test.

The Null Hypothesis we will be testing is: US voters equally rate the police and journalists.

We will do a two-tailed test since we're interested if US voters might not have equal respect for the police and journalists.

### Conducting the chosen test

```
more_police = sum( usVoters$ftjournal < usVoters$ftpolic)
trials = sum( usVoters$ftjournal < usVoters$ftpolic | usVoters$ftjournal > usVoters$ftpolic)
binom.test(more_police , trials)
```

```
##
## Exact binomial test
```

```
##
## data:  more_police and trials
## number of successes = 1008, number of trials = 1795, p-value =
## 2.004e-07
## alternative hypothesis: true probability of success is not equal to 0.5
## 95 percent confidence interval:
##  0.5382418 0.5846764
## sample estimates:
## probability of success
##           0.5615599
```

The p-value is well under 0.05 (the statistical significance we stated in the introduction). That indicates we can reject the null hypothesis that US voters have equal rating for the police and journalists. We extend this conclusion to mean that there is support in the data for the the idea that US voters respect the police more than they do journalists.

Practical significance: The ratio of US voters in this sample who respect the police more to the total number of respondents rating higher either police or journalists is 0.56. While perhaps not earth shattering, in terms of voting preferences that is a significant result. In close elections sometimes decided by less than a percent, a split in opinion of this kind could play a very important role.

## Question 2: Are Republican voters older or younger than Democratic voters?

### Introducing the topic

**Conceptualize** Here we are trying to understand if there is an age component in the differences between the voters of the two largest US political parties.

**Operationalize** The survey itself asks its respondents many questions surrounding political identity. We have several variables which we can consider using to define political identity. First, there are *pid1d* and *pid1r* which ask “Generally speaking, do you usually think of yourself as a Democrat, a Republican, an Independent, or what?”. Each of these reverse the order of which party is listed first. So we can use the data from *pid1d* and *pid1r* (jointly, since 50% of respondents were asked one and 50% the other). We will disregard WHO they voted for, as voters are known to cross party lines (even in this polarized time). For age, we can simply use *birthyr* with the caveat that we only have a yearly granularity for our data (rather than month and day).

**Exploratory Data Analysis** Variables: *pid1d*, *pid1r*, *birthyr*

Types: Both *pid1d* and *pid1r* are ordinal type variables. However, since we are limiting the sample to Republicans(0) and Democrats(1), we can make our new variable “Party” nominal. First, let's determine if we have enough Democrats and Republicans. Using only responses 1 or 2 for *pid1d* and *pid1r* will remove non-responses, Independents, and “others” from the analysis. We will also have to confirm these individuals are voters in 2018 using *voted2018*.

Number of Entries:

```
Age_and_Party <- select(A, c("pid1d", "pid1r",
                             "voted2018", "birthyr",
                             "weight", "weight_spss", "caseid"))

Age_and_Party$Party <-
  ifelse((Age_and_Party$pid1d==2 | Age_and_Party$pid1r == 2), 0,
  ifelse((Age_and_Party$pid1d==1 | Age_and_Party$pid1r == 1), 1, -1))

Age_and_Party <- filter(Age_and_Party, Party >= 0)
Age_and_Party <- filter(Age_and_Party, voted2018 == 1)
paste("Number of Republicans: ", sum(Age_and_Party$Party==0),";",
      "Number of Democrats:   ", sum(Age_and_Party$Party==1))
```

```
## [1] "Number of Republicans: 503 ; Number of Democrats: 725"
```

There are no Null Values for *birthyr*

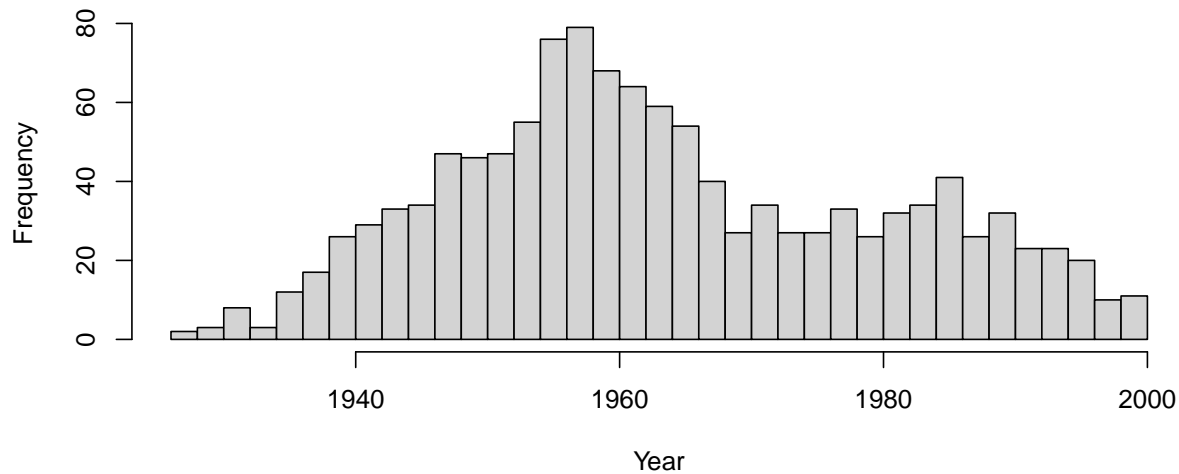
```
paste("Number of Null Values for birthyr:", sum(is.na(Age_and_Party$birthyr)))
```

```
## [1] "Number of Null Values for birthyr: 0"
```

We have 1228 people who claimed they are part of either party, Democrats and Republicans. Now we need to see if the sample distributions are reasonable enough to use CLT. We will draw a histogram of *birthyr* to observe the distribution and any outliers.

```
hist(Age_and_Party$birthyr, breaks = 50,
     main = "Histogram of Democrats and Republicans Birth Years",
     xlab = "Year")
```

### Histogram of Democrats and Republicans Birth Years



The histogram of ages has a somewhat bimodal distribution however, with 1228 responses, we can rely on CLT.

**The Test** Since our sample's political id variable is metric (0 or 1) and *birthyr* is metric, i.i.d (see general assumptions on page 1) and the data cannot be paired, we would argue for using an unpaired two sample t-test. The assumptions using this test is that the variables represented are metric and the data is drawn i.i.d. (See introduction). The *Null Hypothesis* is that the difference ages between Democrats and Republicans is equal to 0.

$$\mu_{D\_age} - \mu_{R\_age} = 0$$

```
t.test(Age_and_Party$birthyr~Age_and_Party$Party)
```

```
##
##  Welch Two Sample t-test
##
## data:  Age_and_Party$birthyr by Age_and_Party$Party
## t = -3.6785, df = 1104.4, p-value = 0.0002459
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -5.235344 -1.593043
## sample estimates:
## mean in group 0 mean in group 1
##      1962.044      1965.458
```

Here we have a statistically significant test with the p-value being less than 0.05. We can reject the null hypothesis that there is no age difference between Democratic and Republican voters. Republican voters are, on average, 3 years older than Democratic voters since they have a mean of a birth year of 1962 and the Democrat voters have a birth year mean of 1965.



```
library(survey)
```

## The Weighted Test

```
## Loading required package: grid

## Loading required package: Matrix

##
## Attaching package: 'Matrix'

## The following objects are masked from 'package:tidyr':
##
##     expand, pack, unpack

## Loading required package: survival

##
## Attaching package: 'survey'

## The following object is masked from 'package:graphics':
##
##     dotchart
```

```
Age_and_Party_Weighted <-
  svydesign(id      = ~caseid,
            weights = ~weight,
            data     = Age_and_Party)

paste("Does Weighting make a difference in the mean birth year?")
```

```
## [1] "Does Weighting make a difference in the mean birth year?"
```

```
paste("Average Birth Year before Weights: ", mean(Age_and_Party$birthyr))
```

```
## [1] "Average Birth Year before Weights: 1964.05944625407"
```

```
paste("Average Birth Year after Weights: ", svymean(~birthyr, Age_and_Party_Weighted))
```

```
## [1] "Average Birth Year after Weights: 1967.45390343812"
```

```
svyttest(birthyr~Party, Age_and_Party_Weighted)
```

```
##
## Design-based t-test
##
## data:  birthyr ~ Party
```

```
## t = 2.8754, df = 1226, p-value = 0.004104
## alternative hypothesis: true difference in mean is not equal to 0
## 95 percent confidence interval:
##  1.243632 6.568711
## sample estimates:
## difference in mean
##           3.906171
```

When applying the weights, the mean birth year of the respondents increased from 1964 to 1967, making the average population of 2018 Democrats and Republican voters younger by nearly 3 years. Even weighted, we still have a statistically significant test is the p-value being less than 0.05. We can again reject the null hypothesis that there is no age difference between Democratic and Republican voters. Using weights, Republican voters are, on average, 3.9 (not 3.4) years older than Democrat voters. It didn't change our outcome of the test, but it slightly modified our practical significance, in which the difference in mean was closer to 4 years than to 3 years.

Why did the weights make a difference? The survey designers added in the weights to make sure the sample is a more correct representation of the US voting population. If they couldn't fully complete their stratified random sampling, or there was a lot of non-respondents, groups could be underrepresented or overrepresented in the sample. Weights are then calculated to give respondents from these groups more or less "power" to the data represented so that one can infer more accurately from the data to the population ... or in other words calibrate the sample to the population.

### Question 3: Do a majority of independent voters believe that the federal investigations of Russian election interference are baseless?

Introduce your topic briefly. (5 points)

*Operationalization* The three most relevant fields look to be `russia16` (Do you think the Russian government probably interfered in the 2016 presidential election to try to help Donald Trump win, or do you think this probably did not happen?), `muellerinv` (Do you approve, disapprove, or neither approve nor disapprove of Robert Mueller’s investigation of Russian interference in the 2016 election?) and `coord16` (Do you think Donald Trump’s 2016 campaign probably coordinated with the Russians, or do you think his campaign probably did not do this?).

The key question is what does it mean “to believe that the federal investigations of Russian election interference are baseless”? Our interpretation is that this phrasing is the equivalent of “to believe that no Russian election interference happened”. Using a measure of approval of the specific Mueller investigation doesn’t seem up to the task, as it is quite possible to “not believe that the investigation is baseless”, but disapprove of Robert Mueller’s work for some other reason.

The remaining variable, measuring opinions whether the Trump campaign coordinated with the Russians, is also tangential to the main question, which refers to the investigation as one of Russian interference, not as an investigation of the Trump campaign’s coordination.

With the exception of possible “No Answers”, the variable `russia16` is a binary variable (yes/no) and therefore metric. If we get some “No Answers” we’ll have to consider eliminating them.

*Population:* We can take the self-defined Independent voters for whom `pid7x` is 3, 4, or 5.

*Gaps:* The main issue here is if we can extend a measure of disbelief in Russian interference as an indicator of belief in the baselessness of the investigation. At the same time, it seems reasonable that if independent voters believed in the existence of Russian interference, they would not find an investigation of such interference baseless. So the variable seems like a good measure.

Perform an exploratory data analysis (EDA) of the relevant variables. (5 points)

Let’s pull out our sample.

```
ind_voters <- A %>%
  filter(pid7x == 3 | pid7x == 4 | pid7x == 5)
ind_voters <- filter(ind_voters, voted2018 == 1)
table(ind_voters$russia16)
```

```
##
##    1    2
## 336 265
```

So the subset of independent voters has no -7 (No Answer) values and is properly binary.

Based on your EDA, select an appropriate hypothesis test. (5 points)

The variable is metric and binary, but that is not a major problem for the CLT. In addition, there is a large sample size. We can run a one-sample t.test.

Our null hypothesis is that exactly 50% of independent voters find the investigation baseless. That translates to a mean of the variable of 1.5, but for more clarity we can just recode the negative responses to 0 and the positive responses to 1.

Since we don't have a clear view to the directionality of a possible rejection of the null hypothesis and since we would be interested in a statistically significant result in either direction, we'll run a two-tailed t-test.

**Conduct your test. (5 points)**

```
test_var <- ifelse(ind_voters$ruusia16 == 1, 1, 0)
t.test(test_var, mu = 0.5)

##
## One Sample t-test
##
## data: test_var
## t = 2.9141, df = 600, p-value = 0.0037
## alternative hypothesis: true mean is not equal to 0.5
## 95 percent confidence interval:
## 0.5192605 0.5988760
## sample estimates:
## mean of x
## 0.5590682
```

Statistically significant result. We reject the null hypothesis. There is evidence in the data in support of the alternative hypothesis. Since we ran a two-tailed test we are able to report that the calculated t-statistic of 2.9141 is in the upper tail, that is pointing to a majority of independent voters believing that Russian interference did in fact occur, and hence by the logic laid out in our introduction, that only a minority of independent voters find the federal investigation baseless.

In terms of practical significance, given the polarization and relative parity of the Democratic and Republican voting blocks, if indeed in the population 56% of independent voters believe that the federal investigation is warranted and that motivates them to vote against the party that benefited from the alleged interference, that could be a very significant effect. In close elections, which are usually zero-sum games, a 6% percent swing from parity to one side is actually a 12% swing in the overall result. In the political campaign world that is massive, even when it is only among a subsection of the electorate. And in this case that happens to be the most sought after group of voters.

## Question 4: Was anger or fear more effective at driving increases in voter turnout from 2016 to 2018?

*Introduction:* This question is complex and touches on several different ideas than need to be unpacked. First we need to find in the dataset some measure of voter fear and anger. Then we need to address if it is at all feasible to reason about “effectiveness in driving voter increases”. Finally, we need to somehow define what voter turnout increases between 2016 and 2018 might be and what sample from the dataset might reflect that.

*Variables* The leading candidates for looking at voter anger and fear are `geangry` and `geafraid` - two ordinal variables measuring the reaction to possible responses to the question “Generally speaking, how do you feel about the way things are going in the country these days?” The scale ranges from 1 (Not at all) to 5 (Extremely).

There are two other candidate survey questions featuring fear and anger, but they address some narrower issues (Donald Trump’s behavior and immigration to the US) and were also randomized (one half of respondents saw one question, the rest saw the other), they would require a lot more assumptions to operationalize.

*Gaps* The crucial gap here is that the question is implicitly one about causality. Further more, it presumes that anger and fear drive (cause) higher turnout and asks which of the causal effects is stronger. We cannot answer that question with this data. What we can look into is whether there’s a difference between the anger/fear responses of those who voted in 2018 after sitting out the 2016 election and the responses of those who didn’t vote in either.

The second gap is closely connected to choosing our sample. The 2018 election did indeed see a historically high turnout, but still that was lower than the turnout of the 2016 election. The reason is that 2016 was a presidential election year, while 2018 was only a midterm election and those typically have lower turnouts than presidential elections.

*Operationalization and population* In order to address the question, while at the same time acknowledging the two significant gaps, we are going to make the case that by comparing a sample of people who voted in 2018, but not in 2016 to a sample of people who didn’t vote in either election, we are indeed looking at “an increase” in turnout. To put it slightly differently, it can be reasonably expected that close to an entirety of the population of people who care about elections and do vote would definitely turn up for the most consequential of elections - that for President. Under that scenario, in your run-of-the-mill regular midterm election, midterm voters can reasonably be expected to be a subset of those who usually vote for President. So if we see a substantial group of voters voting in a midterm election after not voting for President, we can reason that those voters represent an increase in turnout.

For the operationalization we will create a new binary variable `angrier` which will take on a value of 1 if `geangry` is more than `geafraid`, and 0 if it is less. Observations where they are equal will not contribute to the test.

For extracting the two samples, we first start with our `voter_sample` dataframe and use the field `turnout16` (“In 2016, the major candidates for president were Donald Trump for the Republicans and Hillary Clinton for the Democrats. In that election, did you definitely vote, definitely not vote, or are you not completely sure whether you voted?”), `turnout16b` ([IF turnout16=3] “Do you think you probably voted or probably did not vote?”), and `birthyr` to determine our sample of respondents who were eligible in 2016, didn’t vote in 2016 and voted in 2018.

```
voter_sample <- select(A, c("caseid", "birthyr",
                           "voted2018", "turnout16",
                           "turnout16b", "geangry", "geafraid"))
voter_sample$angrier <- ifelse(voter_sample$geangry > voter_sample$geafraid, 1,
                              ifelse(voter_sample$geangry < voter_sample$geafraid, 0, NA))
to_increase <- voter_sample[which(voter_sample$voted2018 == 1 &
                                   (voter_sample$turnout16 == 2 |
```

```
voter_sample$turnout16b == 2) &
voter_sample$birthyr < 1999 &
!is.na(voter_sample$angrier)),]
```

Then for the sample of non-voters we go back to the entire dataset, filter out our 2018 voters and then grab just those who were eligible in 2016, but didn't vote.

```
non_vote <- voter_sample[which((A$turnout16 == 2 | A$turnout16b == 2) & voter_sample$voted2018 == 0 & v
```

**Perform an exploratory data analysis (EDA) of the relevant variables. (5 points)**

Let's look at the sizes of our samples.

```
paste(count(to_increase), count(non_vote))
```

```
## [1] "40 273"
```

Let's look at missing answers for `geangry` and `geafraid`

```
table(to_increase$geangry)
```

```
##
##  1  2  3  4  5
##  6 10 13  7  4
```

```
table(to_increase$geafraid)
```

```
##
## -7  1  2  3  4  5
##  2  8 10 11  4  5
```

```
table(non_vote$geangry)
```

```
##
##  1  2  3  4  5
## 57 56 75 50 35
```

```
table(non_vote$geafraid)
```

```
##
## -7  1  2  3  4  5
##  1 49 78 75 43 27
```

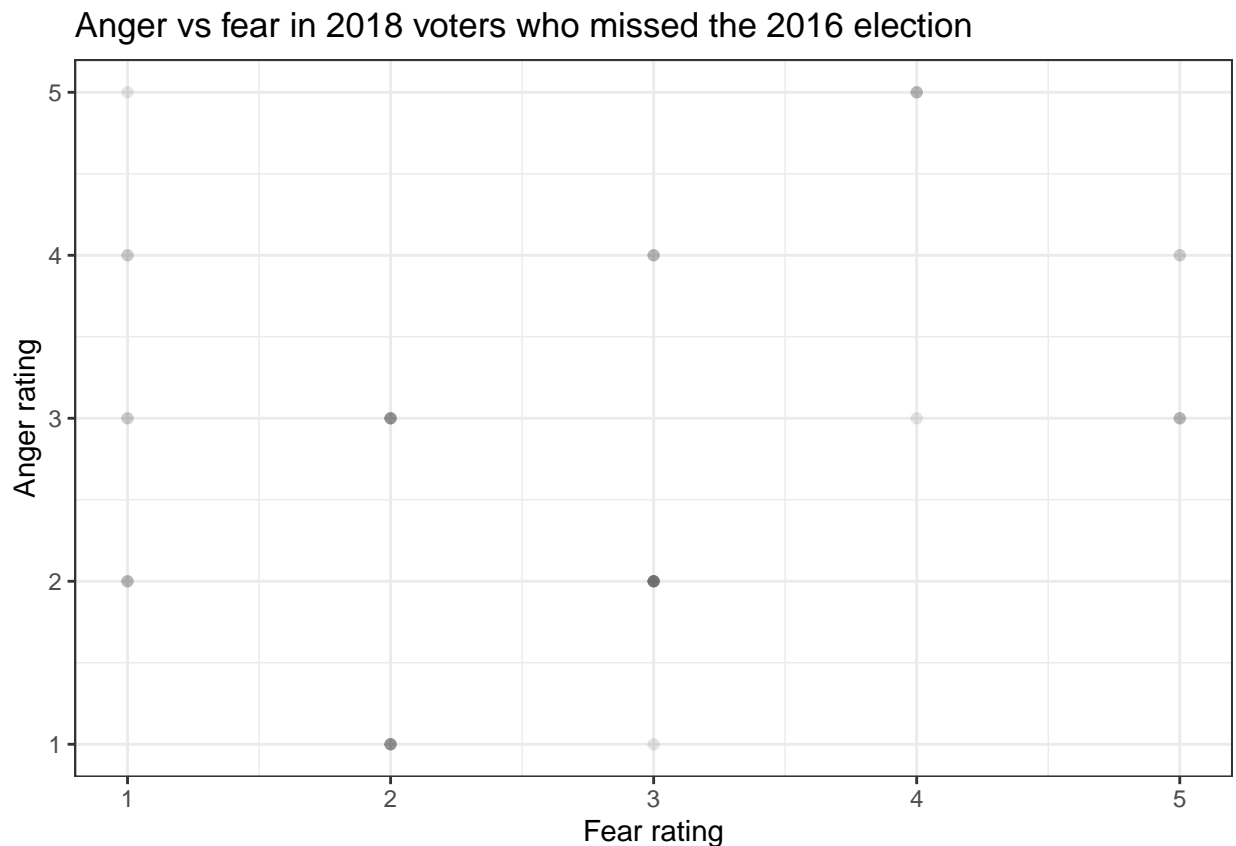
We have 40 observations in one sample and 273 in the other. Of those, only a handful are have value of -7 (No Answer) for `geafraid` or `geangry`. It seems reasonable to omit them.

```
to_increase <- to_increase %>%
  filter(geafraid > 0) %>%
  filter(geangry > 0)

non_vote <- non_vote %>%
  filter(geafraid > 0) %>%
  filter(geangry > 0)
```

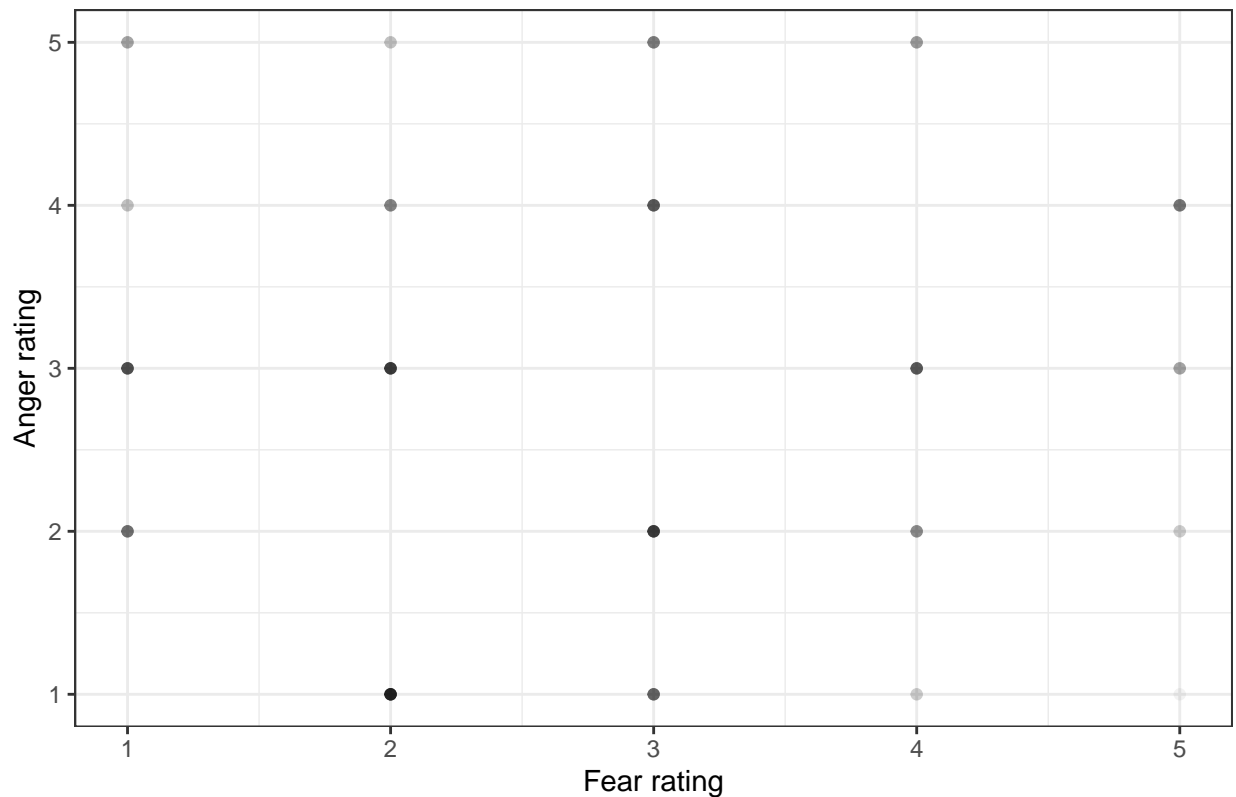
Now we can create two scatter plots to see where the values for `geafraid` and `geangry` fall for the two samples.

```
ggplot(to_increase, aes(x=geafraid, y=geangry) ) +
  geom_point(alpha = 0.1) +
  labs(title = "Anger vs fear in 2018 voters who missed the 2016 election") +
  xlab("Fear rating") +
  ylab("Anger rating") +
  theme_bw()
```



```
ggplot(non_vote, aes(x=geafraid, y=geangry) ) +
  geom_point(alpha = 0.05) +
  labs(title = "Anger vs fear in respondents who didn't vote in 2016 and 2018 elections") +
  xlab("Fear rating") +
  ylab("Anger rating") +
  theme_bw()
```

Anger vs fear in respondents who didn't vote in 2016 and 2018 elections



In the charts above the points below the diagonal represent respondents that are more afraid than angry and the points above - the opposite. They both seem symmetric about the diagonal. We also know that our new variable is binary and since both samples have more than 30 observations we can reasonably rely on the CLT.

Based on your EDA, select an appropriate hypothesis test. (5 points)

We will run a two sample t-test to compare the means of the **angrier** variable for a significance level of 0.05. Newly turned out voters:

```
t.test(to_increase$angrier, non_vote$angrier)

##
##  Welch Two Sample t-test
##
## data:  to_increase$angrier and non_vote$angrier
## t = 0.12587, df = 47.666, p-value = 0.9004
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -0.1651921  0.1872509
## sample estimates:
## mean of x mean of y
## 0.5000000 0.4889706
```

We fail to reject the null hypothesis. Yoda famously said that fear leads to anger, anger leads to hate and hate leads to suffering. Perhaps he had better data and the use of Jedi causal inference. In our case we can



only say that the data do not support the hypothesis that there's a statistically significant difference in the relationship between feelings of fear, on one hand, or feelings of anger, on the other, and increased voter turnout.

## Question 5: Do voters for Hillary Clinton perceive the change in income inequality over the last 20 years differently from how voters for Donald Trump perceive it?

### Introducing our Research Question

The question we are interested in is how Hillary voters view the change in income inequality over the last 20 years versus how Trump voters view the change in income inequality over the last 20 years.

Concept: To evaluate if Hillary voters look at the change in income inequality over the last 20 years differently from the way Trump voters do. We are interested in seeing what people's perception of change in income inequality is in 2018 based on what candidate they preferred in 2016.

Operationlization: The variable `richpoor` captures the answer to the question "Do you think the difference in incomes between rich people and poor people in the United States today is larger, smaller, or the same as it was 20 years ago?". This question was presented to all the participants of the survey. The scale is 1-7 (from a lot larger to a lot smaller):

We will use this variable to understand how the voters viewed the change in income inequality in the last 20 years.

The variable `vote16` captures the answer to the questions "In the 2016 presidential election, who did you vote for? Donald Trump, Hillary Clinton, or someone else?". It takes one of the following three values (along with the representing choices that the participants given to choose from): 1 - Donald Trump 2 - Hillary Clinton 3 - Someone else We will use this to distinguish between Hillary voters and Trump voters.

### Exploratory data analysis (EDA) of the relevant variables

Variables: `richpoor`, `vote16`

Types: `richpoor` is an ordinal variable since it has categories that have an order. Operataions like `>`, `<`, `=` are valid for the 'richpoor' variable. `vote16` is a nominal variable since it has categories that do not have an order. Operataions like `>`, `<`, `=` do not make sense for the `vote16` variable. We are only using the `vote16` variable to separate out the two samples we are interested in comparing.

Number of Entries

```
hillary_voters <- A %>% filter(A$vote16==2)
trump_voters <- A %>% filter(A$vote16==1)

paste("Number of voters who voted for Hillary Clinton in 2016:", length(hillary_voters$vote16))

## [1] "Number of voters who voted for Hillary Clinton in 2016: 898"

paste("Number of voters who voted for Donald Trump in 2016:", length(trump_voters$vote16))

## [1] "Number of voters who voted for Donald Trump in 2016: 770"
```

There are no Null Values for the variable richpoor for Hillary or Trump voters

```
paste("Number of Null Values for richpoor for Hillary voters:", length(which(is.na(hillary_voters$richpoor)
```

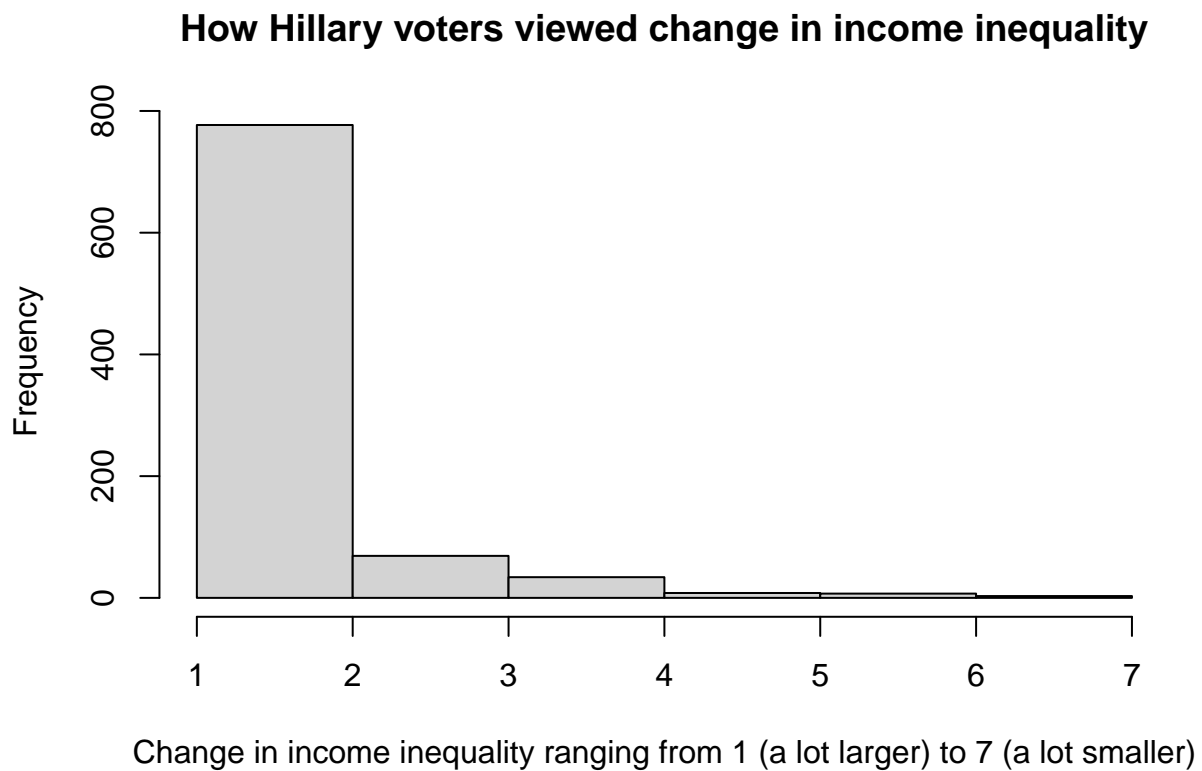
```
## [1] "Number of Null Values for richpoor for Hillary voters: 0"
```

```
paste("Number of Null Values for richpoor for Trump voters:", length(which(is.na(trump_voters$richpoor)
```

```
## [1] "Number of Null Values for richpoor for Trump voters: 0"
```

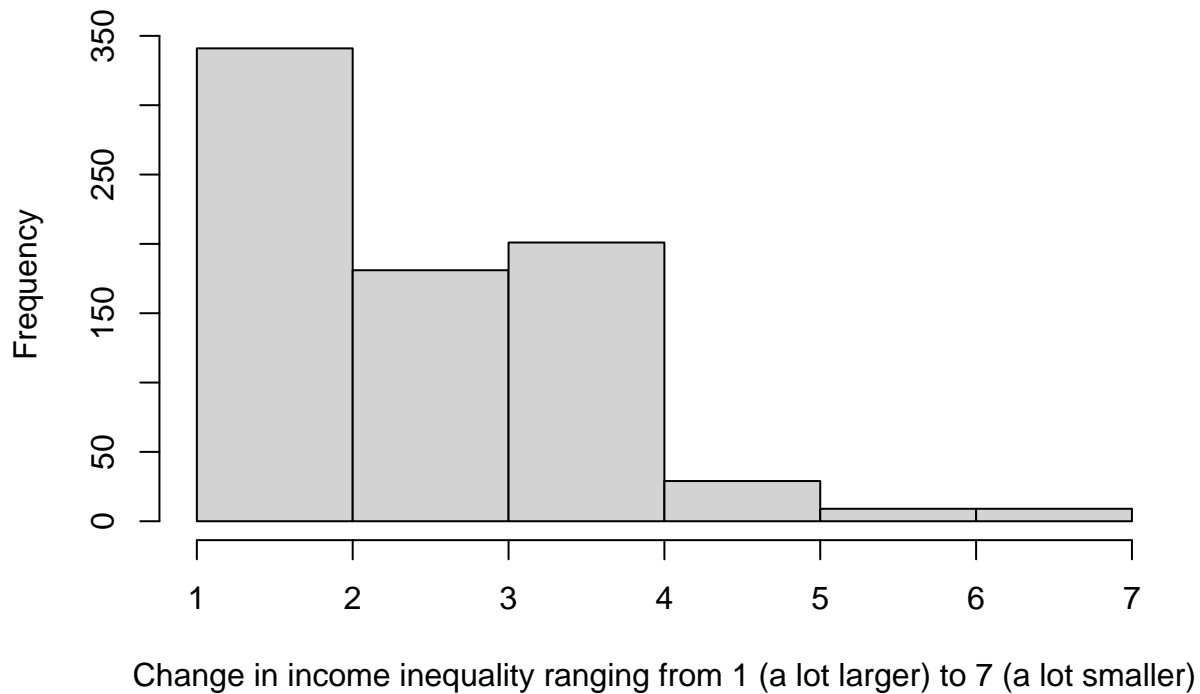
Drawing out the histograms for the richpoor variable for both Hillary and Trump voters to observe sampling distribution.

```
hist(hillary_voters$richpoor, breaks = 7,  
     main = "How Hillary voters viewed change in income inequality",  
     xlab = "Change in income inequality ranging from 1 (a lot larger) to 7 (a lot smaller)")
```



```
hist(trump_voters$richpoor, breaks = 7,  
     main = "How Trump voters viewed change in income inequality",  
     xlab = "Change in income inequality ranging from 1 (a lot larger) to 7 (a lot smaller)")
```

## How Trump voters viewed change in income inequality



Based on our EDA, selecting an appropriate hypothesis test

1. Since we are comparing two samples, it will be a two-sample test.
2. Since the variables are of ordinal scale, parametric tests cannot be used.
3. There is no dependence to rely on and therefore we will use an unpaired test. The distributions from the histograms look slightly skewed towards the left however, since we are intending to use a nonparametric test, the skewness does not affect our choice.

For the above reasons, we narrow our choice of test down to Wilcoxon rank-sum test.

The Null Hypothesis we will be testing is: Hillary voters do not view change in income inequality over the last 20 years differently than Trump voters do.

We will do a two-tailed test since whether there's any difference in the way Hillary voters and Trump voters view change in income inequality over the last 20 years.

Conducting the chosen test

```
wilcox.test(trump_voters$richpoor,hillary_voters$richpoor)
```

```
##  
## Wilcoxon rank sum test with continuity correction
```

```
##
## data: trump_voters$richpoor and hillary_voters$richpoor
## W = 536432, p-value < 2.2e-16
## alternative hypothesis: true location shift is not equal to 0
```

The p-value is less than 0.05 indicating that this test is statistically significant and we can reject the null hypothesis that Hillary voters do not view change in income inequality over the last 20 years any differently than Trump voters do.

## Conclusion

From this test we see that Hillary voters view the change income inequality over the past 20 years differently from how Trump voters view this change.

Practical significance. The total number of pairs is,

```
(sum(A$vote16 == 1 & !is.na(A$richpoor)) * sum(A$vote16 == 2 & !is.na(A$richpoor)))
```

```
## [1] 691460
```

536432 (value of W obtained from our test) out 691460 pairs showed that there was a shift in the perceived rating for change in income inequality between Trump and Hillary voters. The percentage that gets us is

```
paste(536432*100/691460, "%")
```

```
## [1] "77.5796141497701 %"
```

77% is a level at which we feel comfortable saying that the test is practically significant along with being statistically significant.

Common wisdom is that Democratic voters are more concerned about issues of inequality as compared to Republican voters and test results are consistent with this view.