

# fec16 Exploratory Data Analysis

Richard Robbins

March 17, 2022

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Setup</b>	<b>1</b>
<b>3</b>	<b>Troubling Signs</b>	<b>2</b>
3.1	Distribution of General Election Voting Percentages . . . . .	2
3.2	Party Affiliation . . . . .	2
3.3	Multiple Observations For Some Candidates . . . . .	3
<b>4</b>	<b>Mitigation</b>	<b>5</b>
4.1	The Approach . . . . .	5
4.2	Assessing Mitigation Effectiveness . . . . .	5
4.3	Examining “Other Party” Winners . . . . .	6
4.4	The Lone Duplicate Entry Remaining . . . . .	7
4.5	Revisiting that Histogram . . . . .	7
<b>5</b>	<b>What About Homework 10?</b>	<b>7</b>

## 1 Introduction

The w203 Unit 10 homework assignment revolves around the Federal Election Commission 2016 dataset available as an R package. As I worked through that assignment I noticed several anomalies. In the wake of the assignment I decided to go back and take a much closer look at the data. This document summarizes what I found.

## 2 Setup

This review looks primarily at two of the `fec16` datasets, `campaigns` and `results_house`. The working dataset for this exercise, `df`, is formed by an inner join of `campaigns` and `results_house`, after which candidates for whom votes were not recorded are removed. The working dataset is limited to columns of interest. The `candidates` dataset from the `fec2016` collection is touched upon briefly in Section 3.3.

```
campaigns <- fec16::campaigns
candidates <- fec16::candidates
results_house <- fec16::results_house
```

```
df <- inner_join(campaigns, results_house) %>%
  drop_na(general_votes) %>%
  select (cand_id, cand_name, pty_cd, cand_pty_affiliation, party,
          incumbent, ttl_disb, general_votes, general_percent, won, state, district_id)
```

```
## Joining, by = "cand_id"
```

## 3 Troubling Signs

### 3.1 Distribution of General Election Voting Percentages

Figure 1 is a histogram of the general election voting percentages for candidates receiving votes as reflected in the `results_house` dataset. I was surprised to see the asymmetric mass on the left hand side of the histogram. I thought that perhaps seeing a larger group of candidates receiving low percentages might have reflected a third or fourth candidate in a contest between two other more dominant candidates. But I struggled to come up with a satisfactory explanation nonetheless.

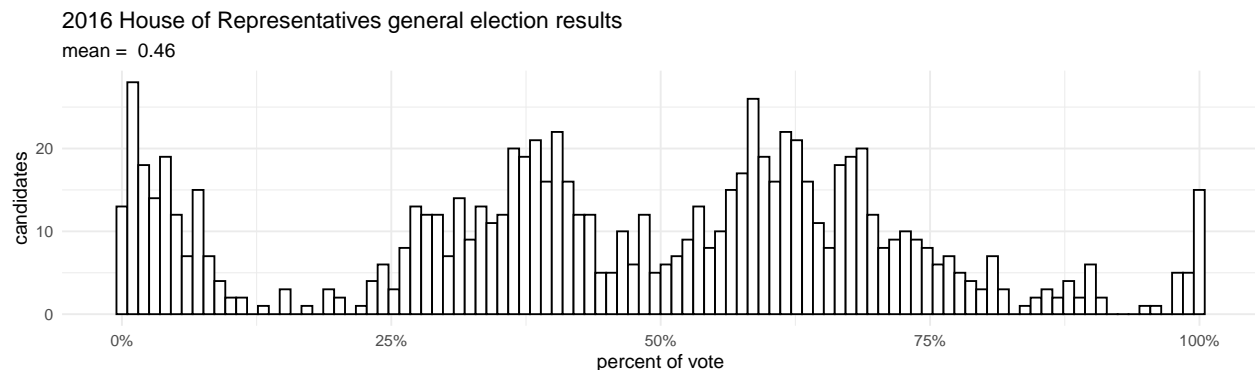


Figure 1: Voting Percentages

### 3.2 Party Affiliation

There are three ways to identify a candidate's party affiliation from the `campaigns` and `results_house` datasets.

1. `df$party` (sourced from `results_house`),
2. `df$cand_pty_affiliation` (sourced from `campaigns`), and
3. `df$pty_cd` (sourced from `campaigns`).

Here are the unique values from each.

```
(unique (df$party))
```

```
## [1] "REP"      "DEM"      "NOP"      "W(R)/R"   "WF"
## [6] "IP"       "LBF"      "NPA"      "CON"      "W"
## [11] "GRE"      "LIB"      "IND"      "R/W"      "U"
## [16] "W(D)/D"   "NLP"      "DFL"      "IDP"      "DNL"
## [21] "WDB"      "MGW"      "NPY"      "WEP"      "CRV"
## [26] "REF"      "R/TRP"    "SID"      "TGP"      "PCC"
## [31] "UPJ"      "N"        "D/IP"     "R/IP"     "D/PRO/WF/IP"
## [36] "R/CON"    "PPD"      "NPP"      "PPT"      "D/R"
```

```
(unique (df$cand_pty_affiliation))
```

```
## [1] "REP" "DEM" "IND" "LIB" "NPA" "CON" "GRE" "UNK" "NNE" "W" "OTH" "IDP"
## [13] "DFL" "NON" "N/A" "UN" "NPP" "PPT"
```

```
(unique (df$pty_cd))
```

```
## [1] 2 1 3
```

If we interpret “REP” and “DEM” as indicating whether the candidate is a Republican or Democrat for the two text fields, and use 1 for Democrat and 2 for Republican for `pty_cd`, we obtain the following.

```
##                Republicans Democrats Other
## party                357         366   157
## cand_pty_affiliation  406         423    51
## pty_cd                406         423    51
```

Using `party` yields 106 more “other party” candidates than the alternatives. The `cand_pty_affiliation` and `pty_cd` fields are equivalent to each other for our purposes. That equivalency holds throughout this exercise.

Let’s see what the numbers look like if we focus only one candidates who were elected. We use `df$won` (sourced from `results_house`) to identify winners..

```
##                Republicans Democrats Other
## party                237         189    64
## cand_pty_affiliation  259         225     6
## pty_cd                259         225     6
```

The numbers are troubling. The dataset indicates that 490 people were elected to the House of Representatives in the 2016 election. However, the number of voting representatives in the House of Representatives is fixed by law at no more than 435. Moreover, after the 2016 general election, all of the voting members of the House of Representatives were either Republicans or Democrats.

### 3.3 Multiple Observations For Some Candidates

While working with the data, I found instances of candidate IDs appearing more than once in `results_house` and, to a lesser degree, a candidate names appearing more than once in the `campaign` dataset.

To explore this, I wrote a pair of functions, one to count the number of times any candidate id (`cand_id`) appears more than once in a dataset, and another to count the number of times a candidate’s name (`cand_name`) appears more than once in a dataset.

```

redundant.names <- function (dataset){
  if ("cand_name" %in% colnames (dataset))
    {dataset %>% group_by (cand_name) %>% count() %>% filter (n>1)}
  else {NA}
}

redundant.ids <- function (dataset){
  if ("cand_id" %in% colnames (dataset))
    {dataset %>% group_by (cand_id) %>% count() %>% filter (n>1)}
  else {NA}
}

```

The table below shows the degree to which the various datasets include more than one observation for a single candidate id or name.

##	dataset	redundant.names	redundant.ids
## 1	campaigns	19	0
## 2	candidates	59	0
## 3	results_house	NA	97
## 4	df	52	52

Lets take a look at the data for a few candidates appearing more than once the working dataset.

```

## # A tibble: 10 x 7
##   cand_id cand_name      cand_pty_affilia~ pty_cd party ttl_disb general_votes
##   <chr>   <chr>         <chr>          <dbl> <chr>   <dbl>         <dbl>
## 1 H0CT030~ DELARUO, ROSA~ DEM          1 DEM    1150879.      192274
## 2 H0CT030~ DELARUO, ROSA~ DEM          1 WF     1150879.       21298
## 3 H0NY290~ REED, THOMAS ~ REP          2 REP    3072934.     136964
## 4 H0NY290~ REED, THOMAS ~ REP          2 CRV    3072934.       16420
## 5 H0NY290~ REED, THOMAS ~ REP          2 REF    3072934.        876
## 6 H0NY290~ REED, THOMAS ~ REP          2 IDP    3072934.       6790
## 7 H2CT021~ COURTNEY, JOS~ DEM          1 DEM    1154847.     186210
## 8 H2CT021~ COURTNEY, JOS~ DEM          1 WF     1154847.       22608
## 9 H2CT051~ ESTY, ELIZABE~ DEM          1 DEM    1447329.     163499
## 10 H2CT051~ ESTY, ELIZABE~ DEM          1 WF     1447329.     15753

```

There is a pattern. The rows for the candidates differ in that it appears that a candidate's total vote is splintered and allocated to several parties.

The campaign spend field (`ttl_disb`) field is consistent for a candidate. The campaign spend information comes from the `campaigns` dataset. From the review above, we know that we do not see redundant entries on candidate ID in that set. So, since we see candidates with multiple entries based on candidate ID in the `results_house` dataset, when we assemble our working dataset by performing an inner join on `campaigns` and `results_house` (which uses the only common field, the candidate ID) we end up replicating campaign spend across multiple rows in the joined dataset. That implies that our working dataset not only has too many other party candidates, the campaign spend is inflated.

I did some research and, at least for the four candidates listed above, I confirmed that the sum of the votes reflected for each above is the total number of votes the candidate received in the election.

We also see that while each of the four candidates listed above is shown as being affiliated with one of the major parties, we see other party affiliations listed. That observation contributes to the inflated number of other party candidates we see in the data.

## 4 Mitigation

### 4.1 The Approach

Based on what I observed and discussed in Section 3.3, in particular, it seems reasonable to group the candidates by ID, party affiliation and, to be prudent, other fields that identify them. Then we derive the vote total for the group by taking the sum of the votes reflected in each row in the group. Party affiliation should be taken from `pty_cd`.

In most cases where there are no duplicate entries, this step will result in the original row for the candidate being preserved. But for candidates with multiple entries based on the candidate ID field, this step will collapse those rows into a single row where the correct number of votes will be reflected.

The cleaned data frame can be derived from the original as follows:

```
df.deduped <- df %>%
  group_by(cand_id, cand_name, cand_pty_affiliation,
           pty_cd, ttl_disb, incumbent, won) %>%
  summarise (general_votes = sum(general_votes))
```

### 4.2 Assessing Mitigation Effectiveness

Let's review what the dataset looks like after the mitigation step described above.

```
(NROW(df))
```

```
## [1] 880
```

```
(NROW(df.deduped))
```

```
## [1] 792
```

```
(redundant.ids(df.deduped))
```

```
## # A tibble: 1 x 2
## # Groups:   cand_id [1]
##   cand_id      n
##   <chr>    <int>
## 1 H2NY03089    2
```

```
(redundant.names(df.deduped))
```

```
## # A tibble: 1 x 2
## # Groups:   cand_name [1]
##   cand_name      n
##   <chr>    <int>
## 1 KING, PETER T. HON.    2
```

1. The number of observations has been reduced from 880 to 792.
2. The number of redundant candidate IDs has been reduced to 1.

3. The number of redundant candidate names has been reduced to 1.

The same candidate is reflected in the second and third items above.

The breakdown by party is as follows:

```
##                Republicans Democrats Other
## cand_pty_affiliation      366      376    50
## pty_cd                   366      376    50
```

And if we isolate on the winners:

```
##                Republicans Democrats Other
## cand_pty_affiliation      236      194     6
## pty_cd                   236      194     6
```

### 4.3 Examining “Other Party” Winners

The six candidates identified as winning but not affiliated with a major party are:

```
## # A tibble: 6 x 8
## # Groups:   cand_id, cand_name, cand_pty_affiliation, pty_cd, ttl_disb,
## #   incumbent [6]
##   cand_id cand_name cand_pty_affili~ pty_cd ttl_disb incumbent won
##   <chr>   <chr>      <chr>          <dbl>   <dbl> <lgl>   <lgl>
## 1 HOMN04~ MCCOLLUM~ DFL              3  966856. TRUE    TRUE
## 2 H2IL13~ DAVIS, R~ UNK              3 2364757. TRUE    TRUE
## 3 H2MN08~ NOLAN, R~ DFL              3 2892902. TRUE    TRUE
## 4 H4MO08~ SMITH, J~ UNK              3 1312777. TRUE    TRUE
## 5 H6PR00~ GONZALEZ~ NPP              3  917851. FALSE   TRUE
## 6 H8MP00~ SABLAN, ~ IND              3   56475. TRUE    TRUE
## # ... with 1 more variable: general_votes <dbl>
```

- Betty McCollum and Richard Nolan are Democrats. The dataset is incorrect.
- Rodney David and Jason Smith are Republicans. The dataset is incorrect.
- Jennifer Gonzalez is neither a Republican nor a Democrat. She represents Puerto Rico as a delegate.
- Geogio Sablan is now a Republican. In 2016 he was neither a Republican nor a Democrat. He represents the Northern Mariana Islands as a delegate.

Because Puerto Rico and the Northern Mariana Islands are territories of the United States, and not states, their representatives in the House of Representatives are delegates with limited voting privileges. Delegates can currently vote in committee and in certain votes on the House floor, but not if their vote would be decisive. Delegates have a marginalized role in Congress and their constituents are not represented in Congress in the same manner as most citizens.

So, while there are 436 people identified as winners in the 2016 House of Representatives election, that include 2 delegates who are not considered voting members of the House. Accordingly, the data is now consistent with the limit on voting members in the House (435) and every voting member has been accounted for as a Republican or Democrat.

## 4.4 The Lone Duplicate Entry Remaining

That leaves one person with duplicate entries in the dataset. That distinction goes to Peter T. King. His record was not merged because when we formed the grouping, to be conservative, we added several indicators from the dataset to prevent merging records that warranted further examination. In this case, the dataset for Mr. King has one row indicating he is not an incumbent who lost and another indicating he is an incumbent who won. In fact, Mr. King was an incumbent and he won re-election in 2016 with 181,506 votes.

```
## # A tibble: 2 x 8
## # Groups:   cand_id, cand_name, cand_pty_affiliation, pty_cd, ttl_disb,
## #   incumbent [2]
##   cand_id cand_name cand_pty_affili~ pty_cd ttl_disb incumbent won
##   <chr>   <chr>      <chr>          <dbl>   <dbl> <lgl>    <lgl>
## 1 H2NY03~ KING, PE~ REP             2 1310730. FALSE    FALSE
## 2 H2NY03~ KING, PE~ REP             2 1310730. TRUE     TRUE
## # ... with 1 more variable: general_votes <dbl>
```

## 4.5 Revisiting that Histogram

Unfortunately reconstructing the histogram presented above will take more effort as all of the vote percentages for the merged duplicated rows would need to be recalculated.

## 5 What About Homework 10?

A review of the model that I built for Unit 10 is beyond what I aim to achieve in this summary. However, I did redo much of my work using the refined dataset that reflects the changes described here.

Had I used this dataset I would have submitted a different model. Using the original dataset, I concluded that there was a statistically significant difference between Republicans and Democrats when it came to the impact of campaign spending on votes received. In my model the constant associated with being a Democrat was higher than that of being a Republican, however, the effectiveness of campaign spend for Republicans outpaced that of Democrats. When I repeat my work using this data, that difference disappears. With this data, I found no difference between Republicans and Democrats. In both cases there was a significant difference between being affiliated with one of the major parties or not. With the cleaner data, the total residuals dropped and adjusted  $R^2$  went from .258 to .426.