

Analyzing the 2016 US Presidential Election

We analyze returns from the 2012 and 2016 elections in order to understand the social and demographic trends that may have contributed to Donald Trump's victory in 2016. We will first examine how Republican vote share at the county level has changed from 2012 to 2016. Then, we will look at four variables that were prominent in the discourse around the election – race, education, unemployment, and immigration – to see how well they predict GOP electoral gains at the county level.

We will be working with three datasets. The first, `election2012.csv`, has one observation per county and contains the following variables:

Name	Description
<code>FIPS</code>	FIPS code (unique county identifier)
<code>state</code>	State abbreviation
<code>county</code>	County name
<code>votes_dem_12</code>	Number of votes cast for Democratic candidate, 2012 election
<code>votes_gop_12</code>	Number of votes cast for Republican candidate, 2012 election
<code>votes_total_12</code>	Total number of votes cast in 2012 election

The second, `election2016.csv`, has the same data structure and similar variable names but reports data for the 2016 presidential election.

The third dataset, `county.csv`, includes social and demographic characteristics for each county:

Name	Description
<code>FIPS</code>	FIPS code (unique county identifier)
<code>pct_for_born15</code>	Percent of county's population that is "foreign born" according to the U.S. Census, meaning anyone who is not a U.S. citizen at birth (measured over 2011-2015)
<code>pct_bach_deg15</code>	Percent of county population holding a Bachelor's degree or above (2011-2015)
<code>pct_non_white15</code>	Percent of county population that is not white (2011-2015)
<code>pct_unemp16</code>	Percent of county population that is unemployed, BLS estimates (average, Jan-Oct 2016)
<code>pct_unemp12</code>	Percent of county population that is unemployed, BLS estimates (average, Jan-Oct 2012)

Question 1

Start by load all three datasets. Merge the three datasets by FIPS code to construct one complete data file for analysis. Check your merge to see how many observations came in from all three sources. Did you lose much data? Finally, perform listwise deletion (hint: check section 3.2 of *QSS*) of missing values on the full dataset. Did you lose much data? Get whatever characteristics you can on the data you lost. What can you say about these observations?

Question 2

Compute the Republican vote share as a proportion of total votes, in 2012 as well as in 2016. Also compute the percent difference in this Republican vote share variable from the 2012 to 2016 election. Plot the distribution of this percent difference, with a red line at the median.

Then, subset your data to just the battleground states: Florida, North Carolina, Ohio, Pennsylvania, New Hampshire, Michigan, Wisconsin, Iowa, Nevada, Colorado, and Virginia. Plot the distribution of the same variable in this sample, with a red line at the sample median.

Question 3

Create a county-level map of the United States, with counties where Democrats got a larger vote share in 2016 than 2012 in blue, and counties where the Republican vote share increased in red. We also want the intensity of the color to depend on the magnitude of the Democratic or Republican gains. To create this map, you will need to take the following steps:

1. If you have not yet done so, install the `maps` package and load it.
2. Take the `county.fips` dataset, which comes with the `maps` library, and perform a merge with the dataset you were working with in the last question, by FIPS code – that is, make sure that your merged data contains all the observations from `county.fips`, in their original order (hint: use the `all.x` argument).
3. As you know, `alpha` values are typically used in the interval $[0,1]$. One way to normalize data to this range is to calculate, for a vector x :

$$\frac{x - \min(x)}{\max(x) - \min(x)}$$

Use this normalization strategy on the Republican vote share variable you calculated above and store the result in an object.

Then, use the `rgb()` function to create a vector of appropriate colors (red for Republican gains, blue for Republican losses), using the vector you created for the `alpha` argument inside that function. Think carefully about what a Republican *gain* and *loss* mean in relation to the vote share variable you calculated. Finally, use this vector of colors within the `map()` function. Include the `lty = 0` option to get rid of black borders around the states.

Comment on your results. In what parts of the U.S. did Republicans make the most significant gains? What other interesting patterns do you observe?

Question 4

Run a regression of percent change in Republican vote share from 2012 to 2016 on percent foreign-born, percent holding a Bachelor's degree or above, percent non-white, and percent unemployed. Interpret your results.

Question 5

We will now see which counties had the most surprising election results in 2016 given our predictions based on the previous election. To do so, first regress 2012 Republican vote share on percent foreign-born, percent with a Bachelor's degree or above, percent non-white, and percent unemployed in 2012. For the first three, you can use the variables ending in 15 since these are the most recent available Census estimates, which are

averaged over the period 2011-15. Then predict 2016 Republican vote share in each county using the same 2011-15 variables and percent unemployed in 2016. Compute the prediction error, which is the predicted Republican vote share subtracted from the observed value in 2016. Create a county-level map with counties colored in red where the observed value was higher than the prediction and blue otherwise. Use double the absolute value of the prediction error as the intensity of the color (the `rgb()` `alpha` parameter). Comment on the results.

Question 6

Subset the data to the counties with the largest overpredictions and underpredictions of Republican vote share based on the last question (take the top and bottom quantiles of prediction error).

Create some histograms using these subsets, with black lines at the medians in the subset and red lines at the medians from the full data. Are the counties that defied our expectations unusual in any interesting ways?