

# The Divided States of America

---

Dmitry Mikhaylov and Keita Miyaki

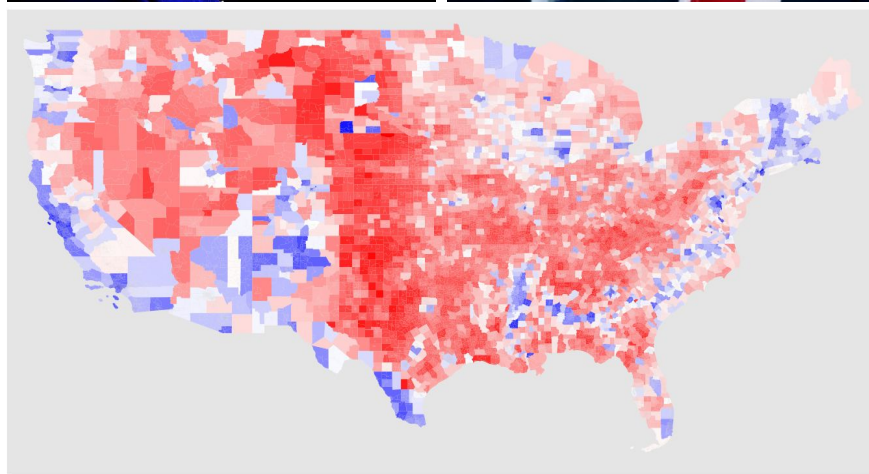
<https://github.com/keita-dc/FlatironProject5>

Hmm...

# Our questions

What did we ask?

1. Can census data tell the winner of 2016 presidential election at a county level?
2. What features matter most?

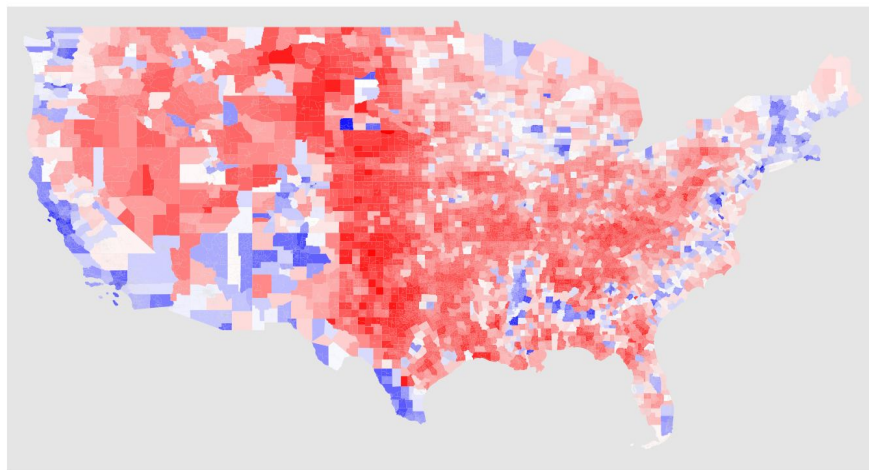


Aha!

# Our discoveries

What did we learn?

1. The nation is so divided that most models can predict the county-level results well
2. Race and Urbanization (housing, transportation, etc) do matter most



# Our Dataset

## US Census

American Community  
Survey 5-Year Data (2017)

- 260 features after cleaning
  - 3,112 countries (excluding Puerto Rico and Alaska)
-

# Our Features

## ACS “Data Profiles”

*Data Profiles contain broad social, economic, housing, and demographic information. The data are presented as population counts and percentages. There are over 1,000 variables in this dataset.*

<https://www.census.gov/data/developers/data-sets/acs-5year.html>

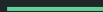
---

# Our Features

## Data Cleaning

### ACS “Data Profiles”

- Only percentage point estimates
- Only features provide county-level data
- Only features with range between 0 and 100



# Our Feature Examples

## Demography

- Race
- Age
- Sex

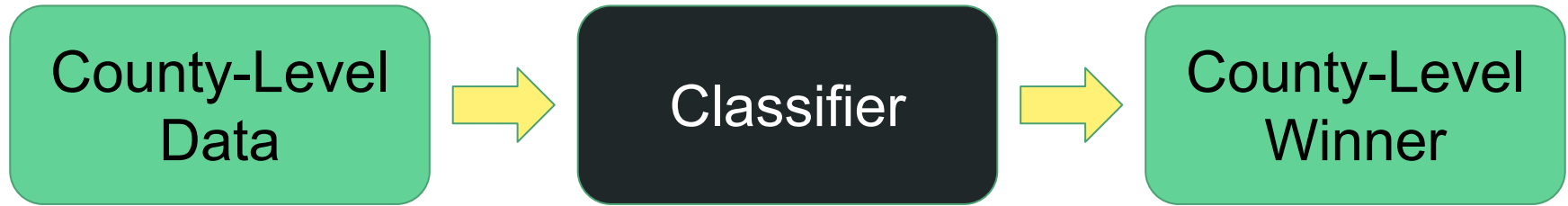
## Urbanization

- Apartments/house
- Means of transportation

## Income

- Income range
-

# Our Analysis



## Census ACS Data

- 260 Features
- 3,112 Counties

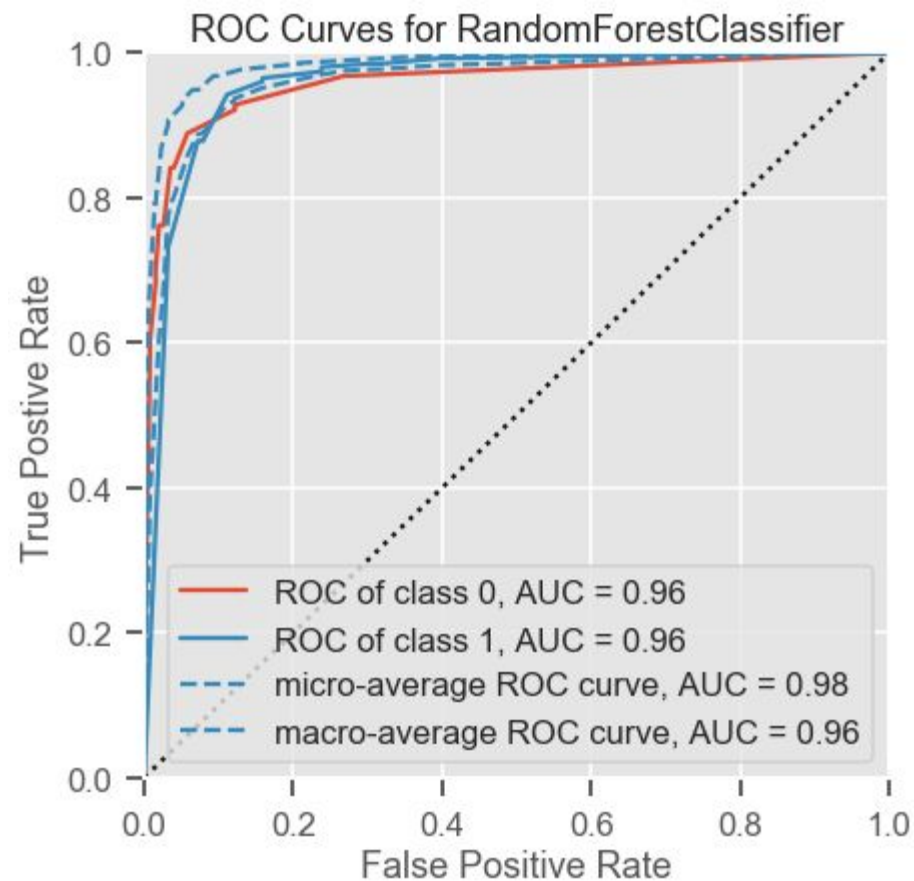
## Models

- Random Forest
- SVM
- XGBoost

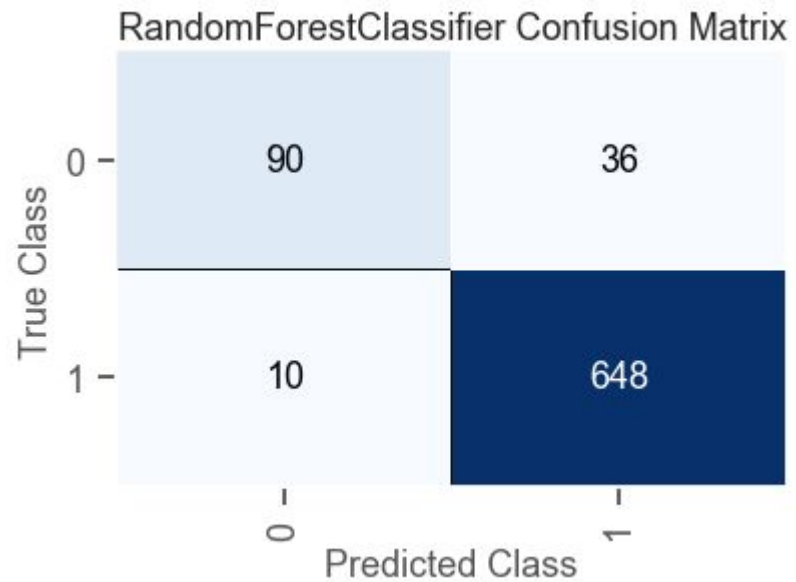
## Candidates

- Hillary ( 0 )
- Trump ( 1 )





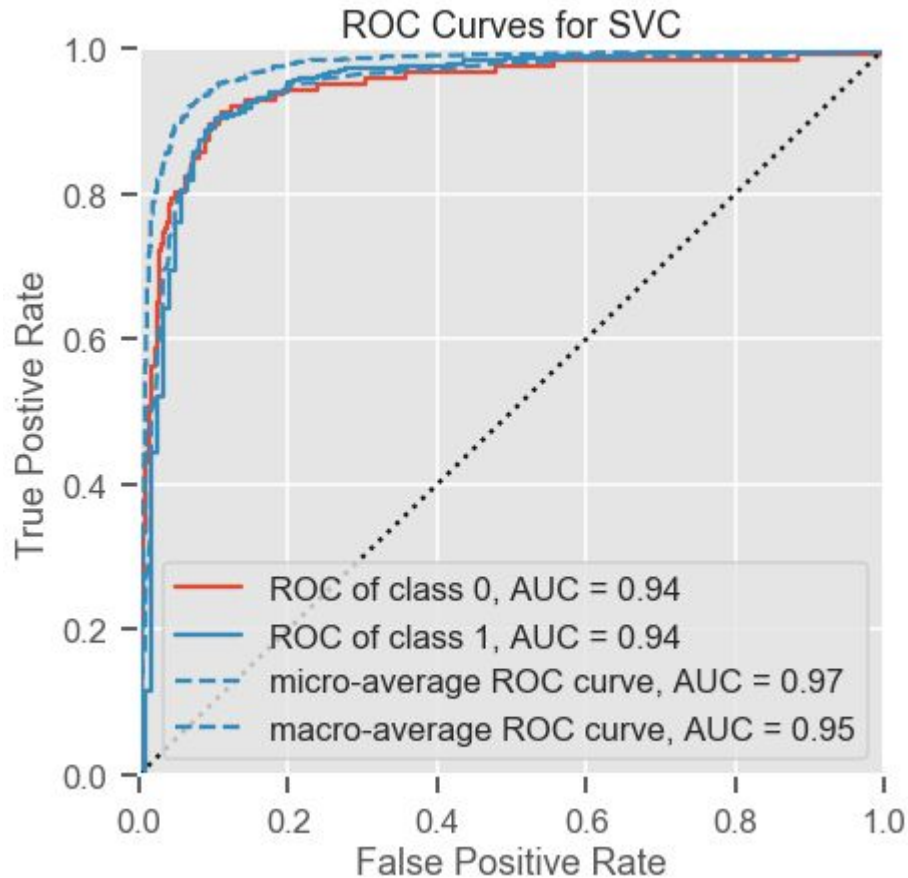
# Random Forest

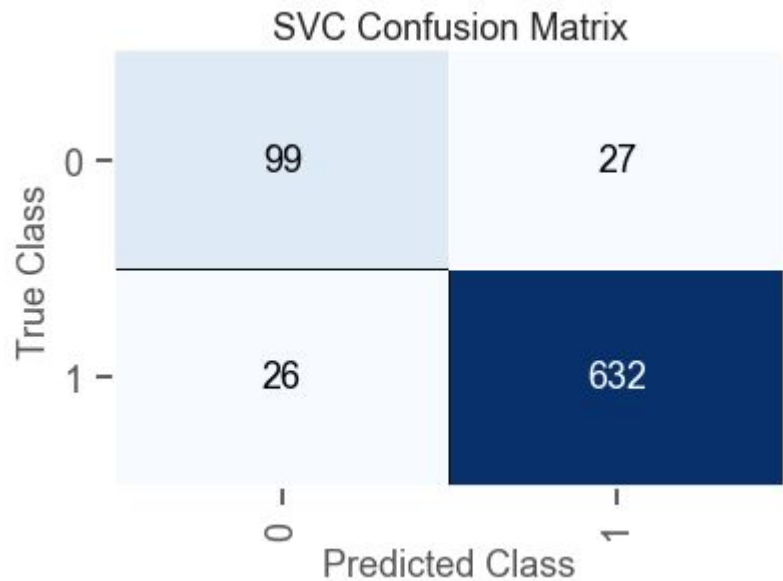


# Random Forest

---

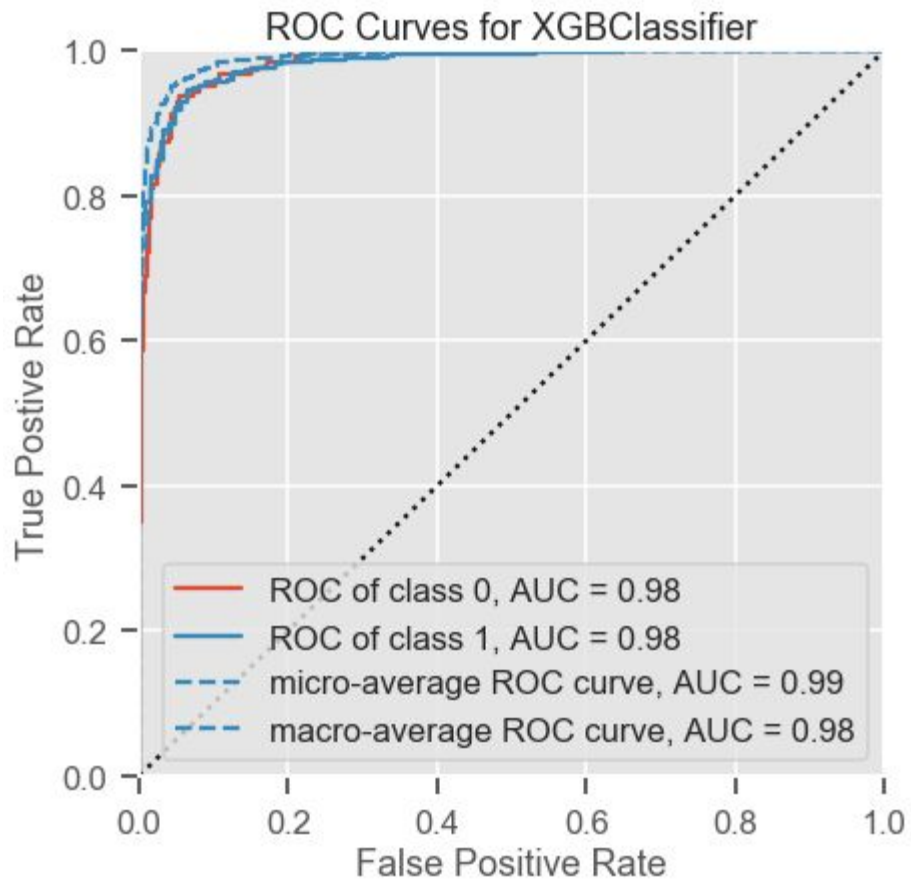
# Support Vector Machine





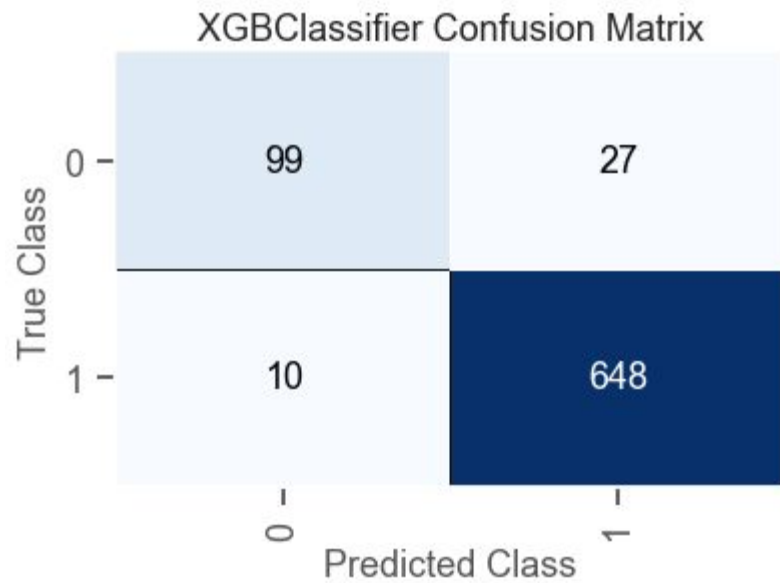
# Support Vector Machine

---

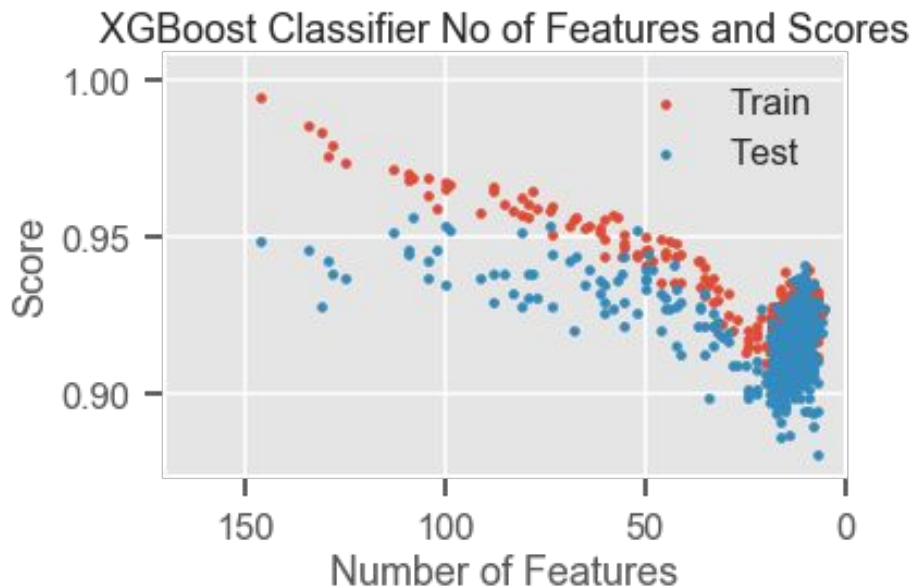


# XGBoost

---



# XGBoost



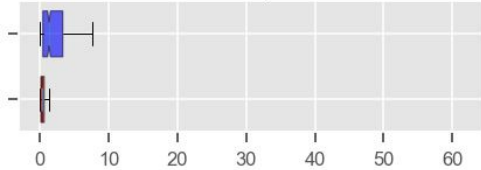
# Feature Selection in XGBoost by Regularization

---

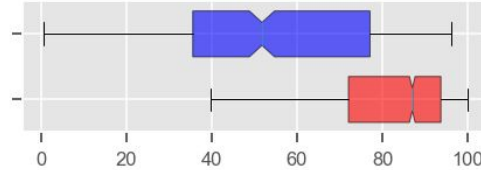
# Important Features

## Regularized XGBoost All Features

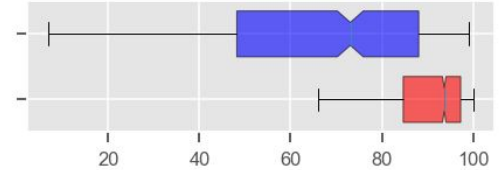
Commuting To Work Workers 16 Years And Over Public Transportation (Excluding Taxicab): 0.29



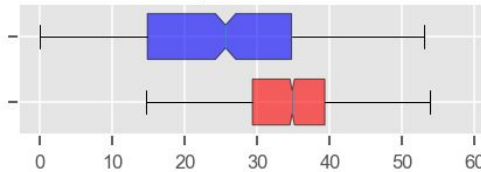
Hispanic Or Latino And Race Total Population Not Hispanic Or Latino White Alone: 0.28



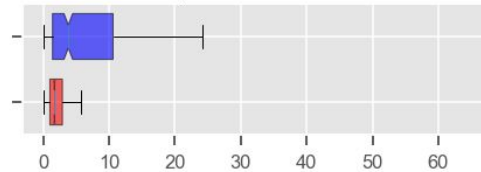
Race Alone Or In Combination With One Or More Other Races Total Population White: 0.18



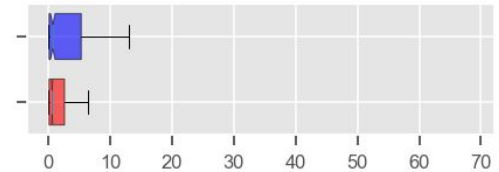
Selected Monthly Owner Costs (Smoc) Housing Units Without A Mortgage \$250 To \$399: 0.10



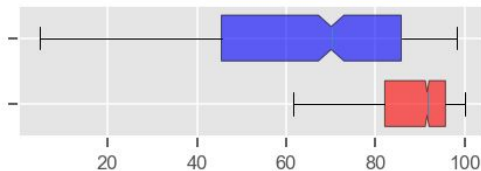
Value Owner-Occupied Units \$500 000 To \$999 999: 0.05



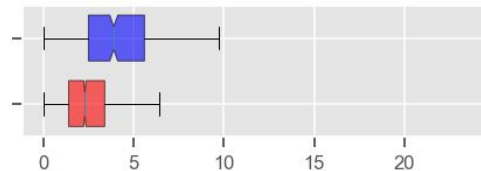
House Heating Fuel Occupied Housing Units Fuel Oil Kerosene Etc.: 0.04



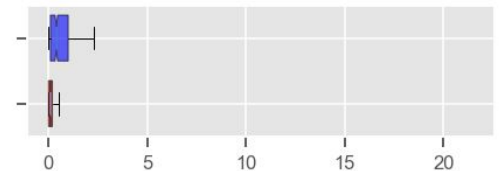
Race Total Population One Race White: 0.04



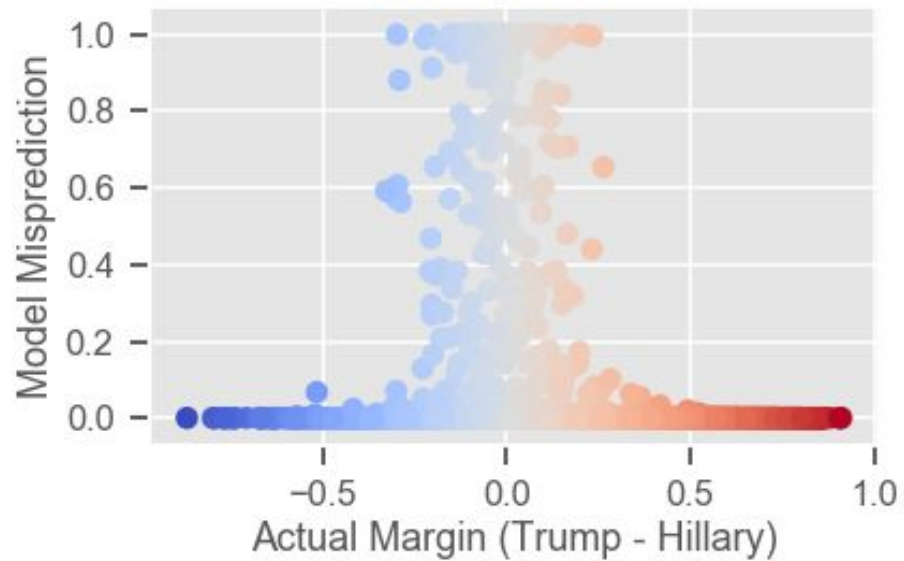
Units In Structure Total Housing Units 3 Or 4 Units: 0.03



Race Total Population One Race Asian Chinese: 0.00



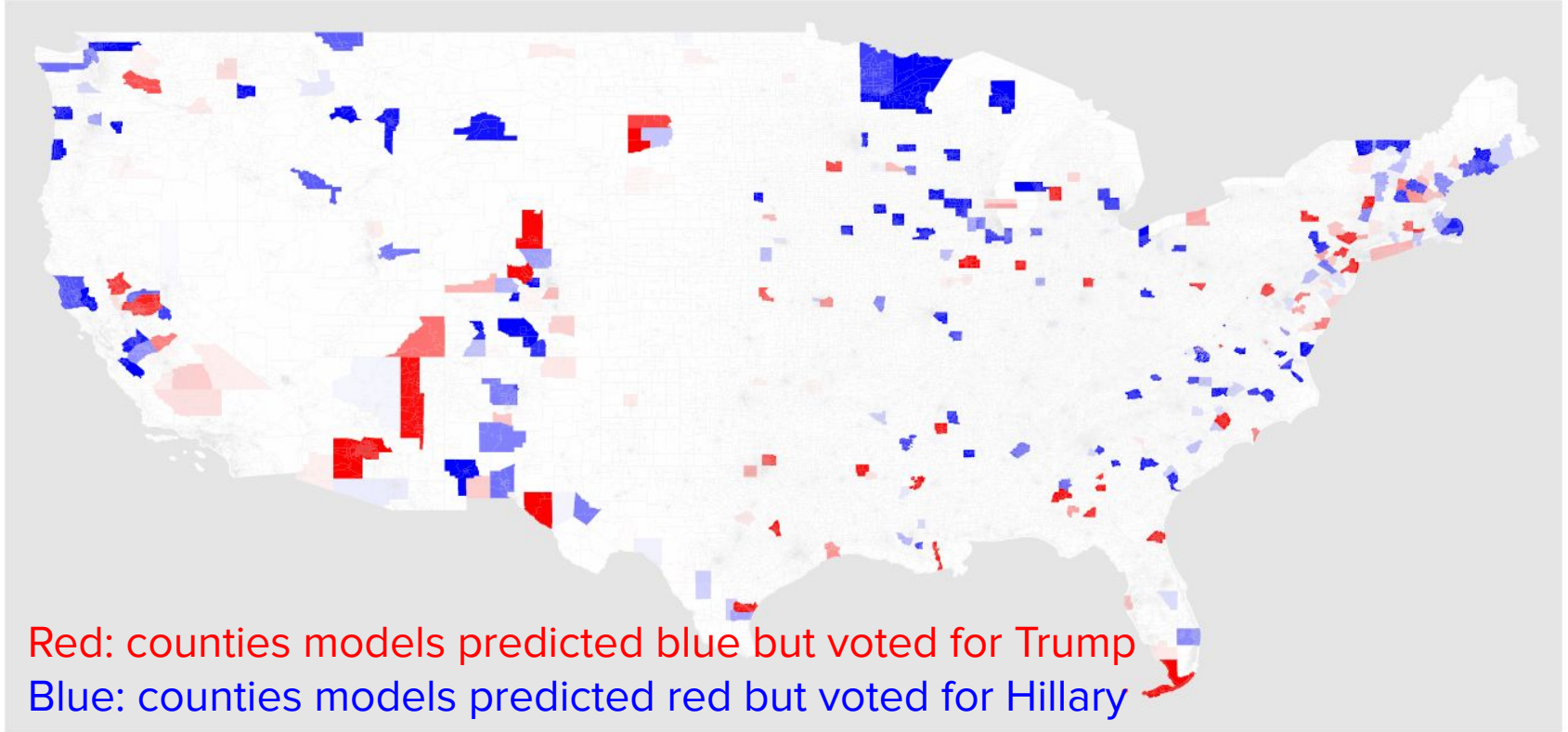




Margin  
vs  
Misprediction

---

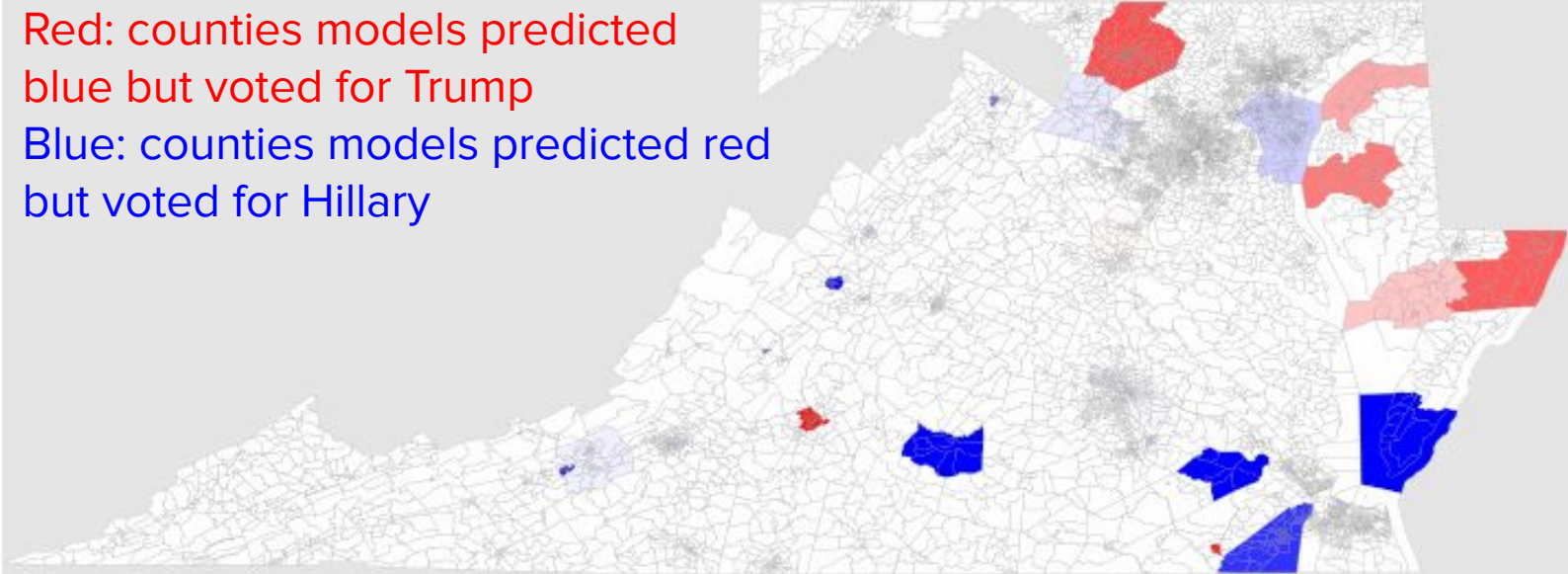
# Confusion Map - Where XGB Mispredicted



# Confusion Map - Where XGB Mispredicted

Red: counties models predicted blue but voted for Trump

Blue: counties models predicted red but voted for Hillary



THANK YOU!  
&  
GO VOTE!

