

A Multivariate Analysis on U.S Presidential Election Results From 1992

Jared Thacker

Objective

1. Is there a difference in census information between counties that Bill Clinton lost and won?
2. Can we predict the results of an election using *only* census information and ignoring time? Baseline model?
3. Which variables are the most important?

Democratic Candidate: Bill Clinton



Incumbent: George H. W. Bush

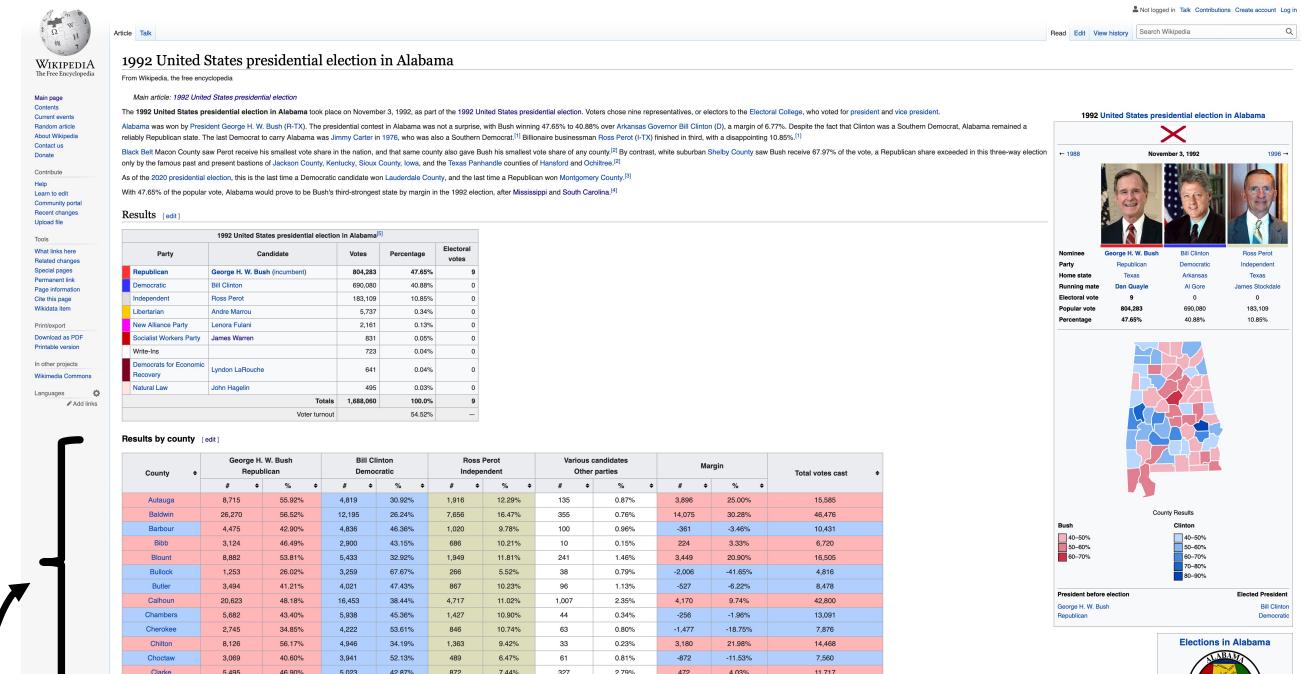
Outline

1. Data-Sourcing
2. The Data
3. Exploratory Data Analysis
4. Question 1/Question 3
5. Question 2
6. Conclusion
7. Extra Graphs
8. Useful R Packages



Data-Sourcing

- ▶ Two Sources
 - ▶ Ufl.edu - census variables by county
 - ▶ Wikipedia - election results by county
- ▶ EXTENSIVE Data-cleaning was needed



An example of a Wikipedia table that was scraped

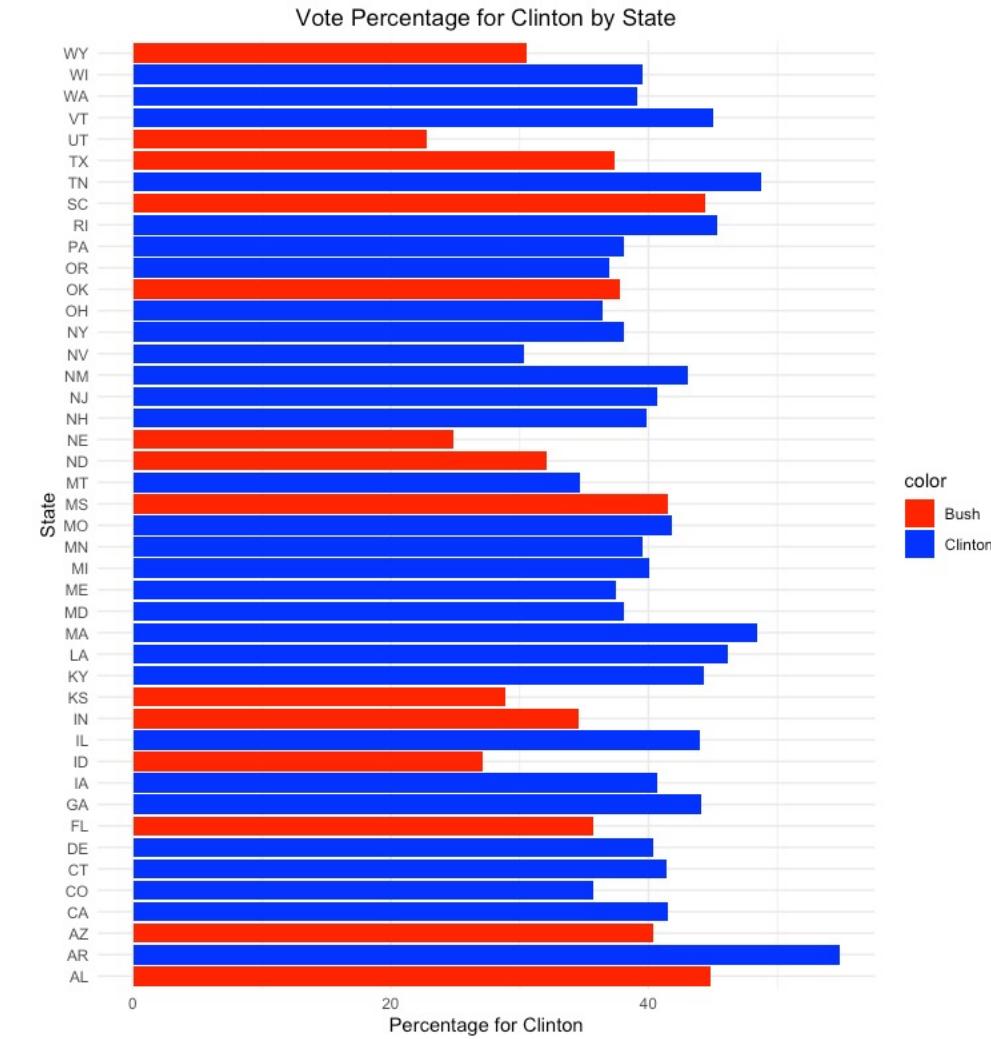
The Data

- ▶ Census Variables - Predictor Variables
 - ▶ Median Age (years) - num./Cont.
 - ▶ Mean Savings (\$) - num./Cont.
 - ▶ Per Capita (PC) Income (\$) - Num/Cont
 - ▶ Percent In Poverty (%) - num./cont.
 - ▶ Percent Veterans (%) - num./cont.
 - ▶ Percent Female (%) - num./cont.
 - ▶ Population Density (?) - num./cont.
 - ▶ Percent In Nursing Homes (%) - num./cont.
 - ▶ Crime Index Per Capita (?)-numerical/cont.?

- ▶ Response
 - ▶ Clinton Win - Binary, nominal
 - 2413 Observations (counties)
 - 9 predictors, 1 response
 - 1220 counties that Bush won
 - 1193 counties that Clinton won
 - Subpopulations almost equivalent

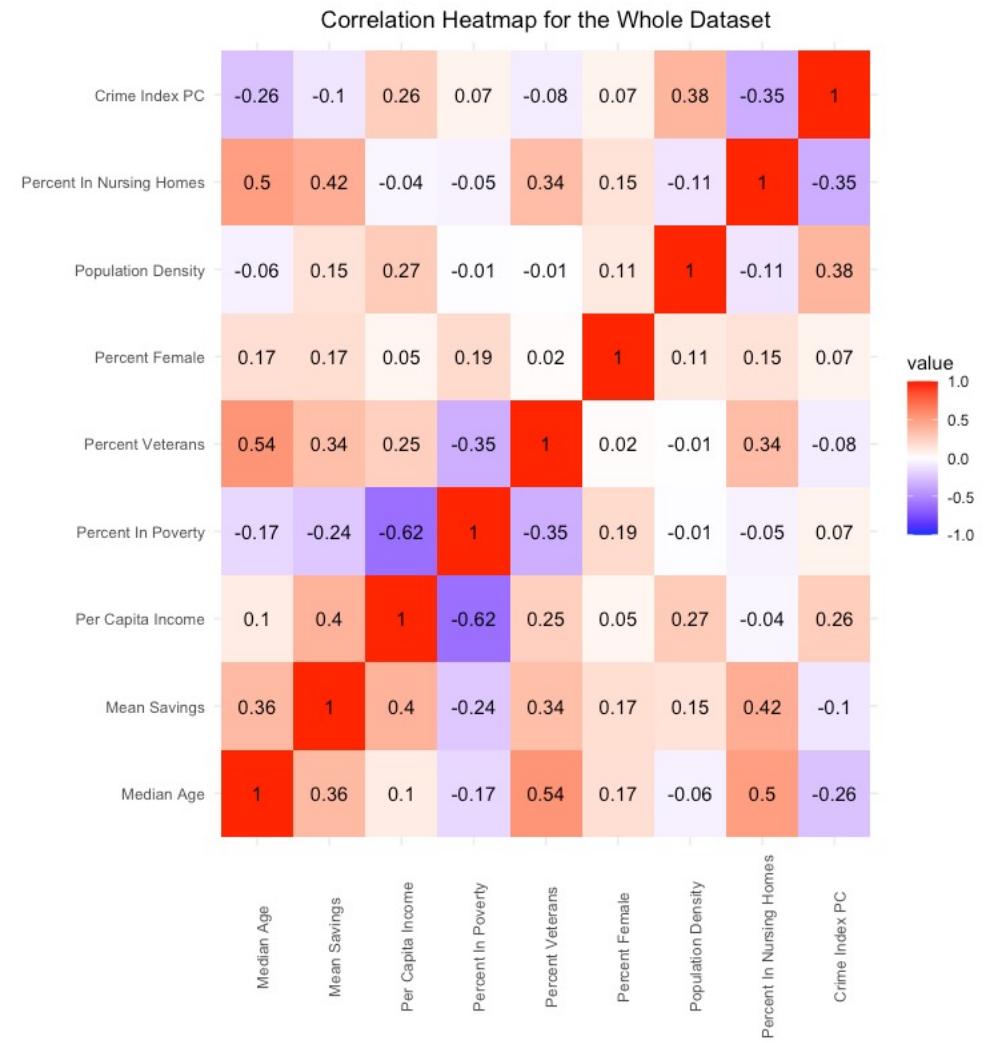
| County Name | ...Census Variables... | Clinton Win |
|-------------|------------------------|-------------|
| Autauga | | FALSE |
| Baldwin | | FALSE |
| Barbour | | TRUE |
| Blount | | FALSE |

Exploratory Data Analysis (EDA)



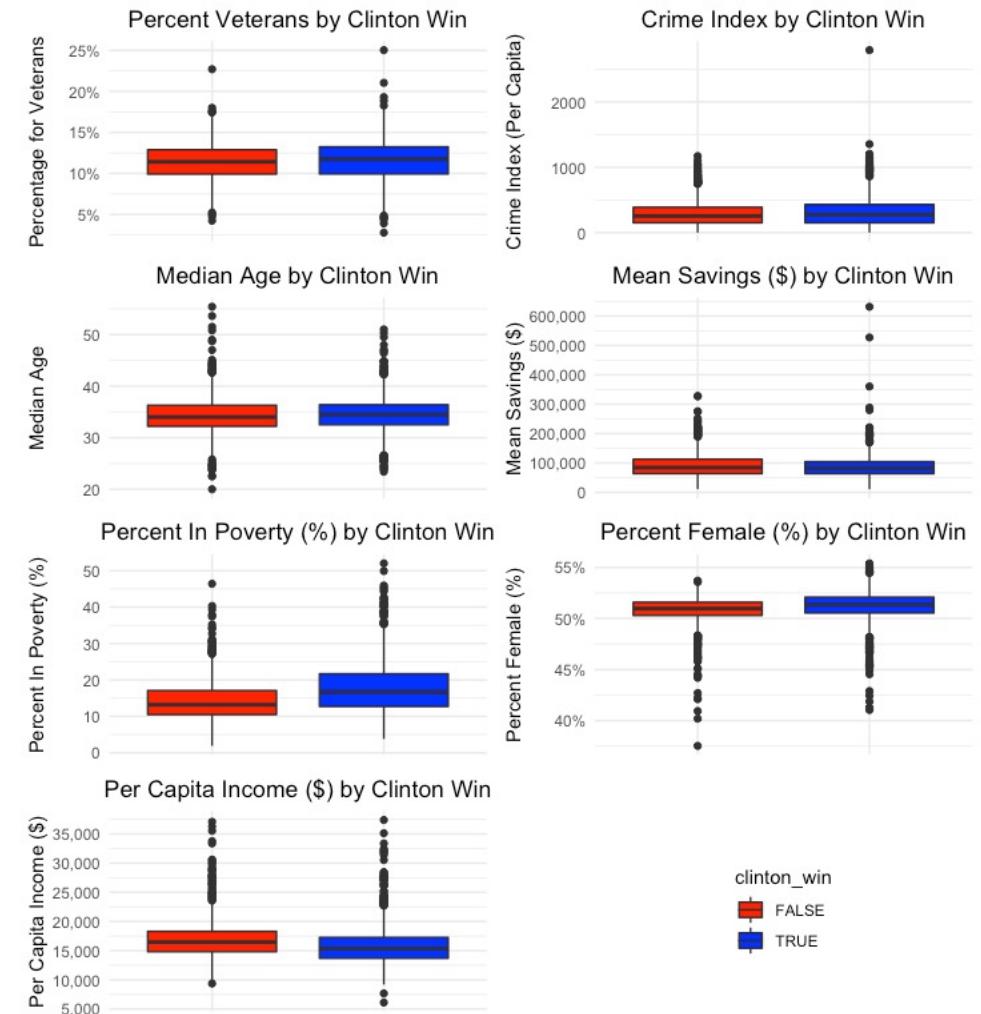
More EDA

- ▶ Moderately strong correlation between variables
- ▶ PCA is appropriate



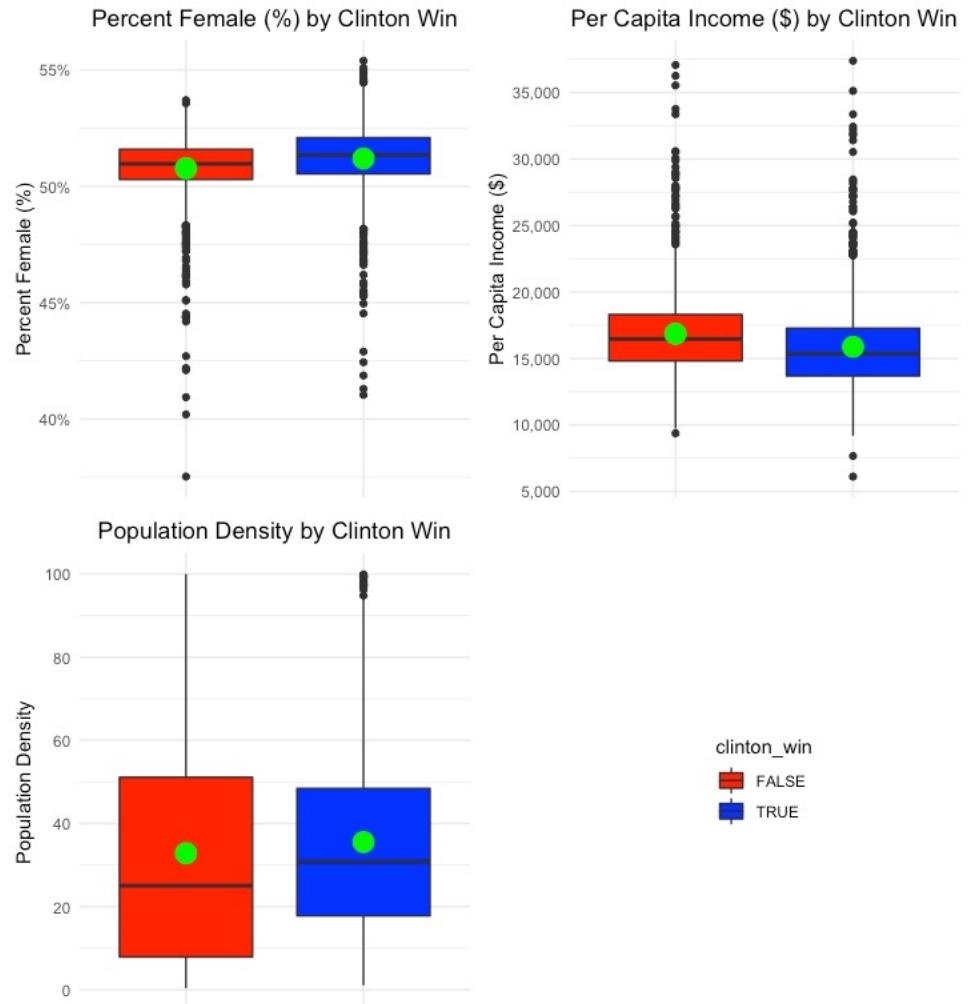
More EDA

- ▶ Appearance: small difference between different census measure
- ▶ Large sample size -> more power



Question 1: MANOVA

- ▶ Use Wilk's test
- ▶ F test statistic: 11.493
- ▶ P -value < 0.00001
- ▶ Important variables (ANOVA)
 - ▶ Mean Savings: $p < 0.001$
 - ▶ PC Income: $p < 0.0001$
 - ▶ % Female: $p < 0.0001$
 - ▶ Population Dens.: $p < 0.005$
 - ▶ % in poverty: $p < 0.0001$
- ▶ Random Forest tells a similar story for variable significance

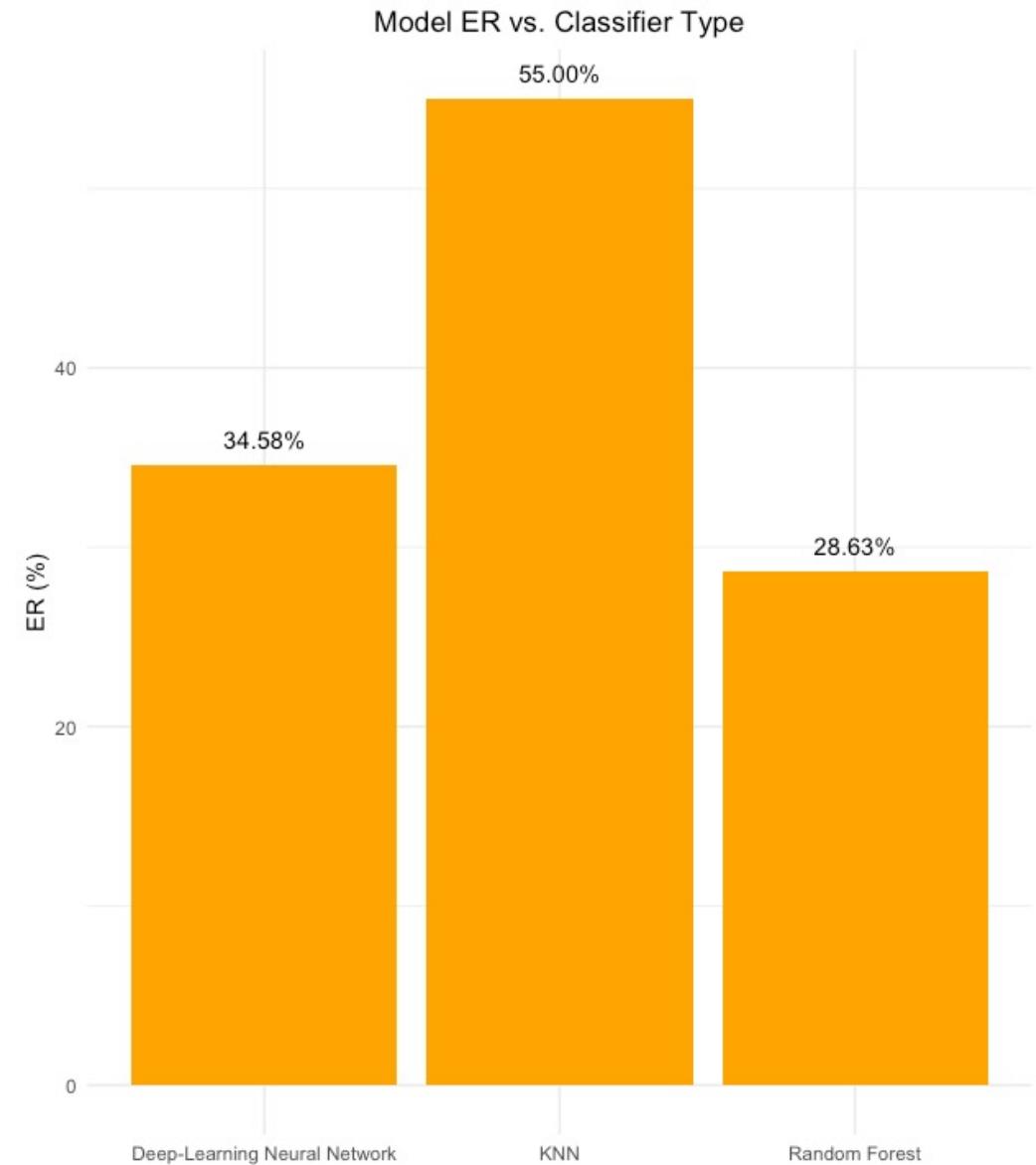


Question 2: Predictive Modeling

- ▶ Three Models
 - ▶ KNN, Random Forest, Dense Deep-Learning Neural Network (DLNN)
 - ▶ Train test split: 80%/20%
- ▶ KNN
 - ▶ K=5
- ▶ Random Forest
 - ▶ # of Trees: 150
 - ▶ variables at each split : 1
- ▶ DLNN (All were similarly bad)
 - ▶ # of hidden layers: 5
 - ▶ # of hidden units:
 - ▶ 32
 - ▶ 64
 - ▶ 128
 - ▶ 64
 - ▶ 32
 - ▶ Dropout rate: 10%
 - ▶ Activation function: Sigmoid
 - ▶ Loss function = “binary cross entropy”
 - ▶ Optimizer = “RMSprop”

Question 2: Predictive Modeling

- ▶ Random forest - highest performer
- ▶ Neural Network
 - ▶ Not enough data (inconclusive)
- ▶ There are more ML models
 - ▶ I chose just three
- ▶ Baseline Model - Random Forest



Question 3: Variable Importance

- ▶ Random forest - highest performer: Population Density
- ▶ Random forest importance is *very similar as* ANOVA results

population_density
percent_in_poverty
per_capita_income
percent_female
percent_veterans
mean_savings
crime_index_PC
percent_in_nursing_homes
median_age

Variable Importance Plot from Random Forest

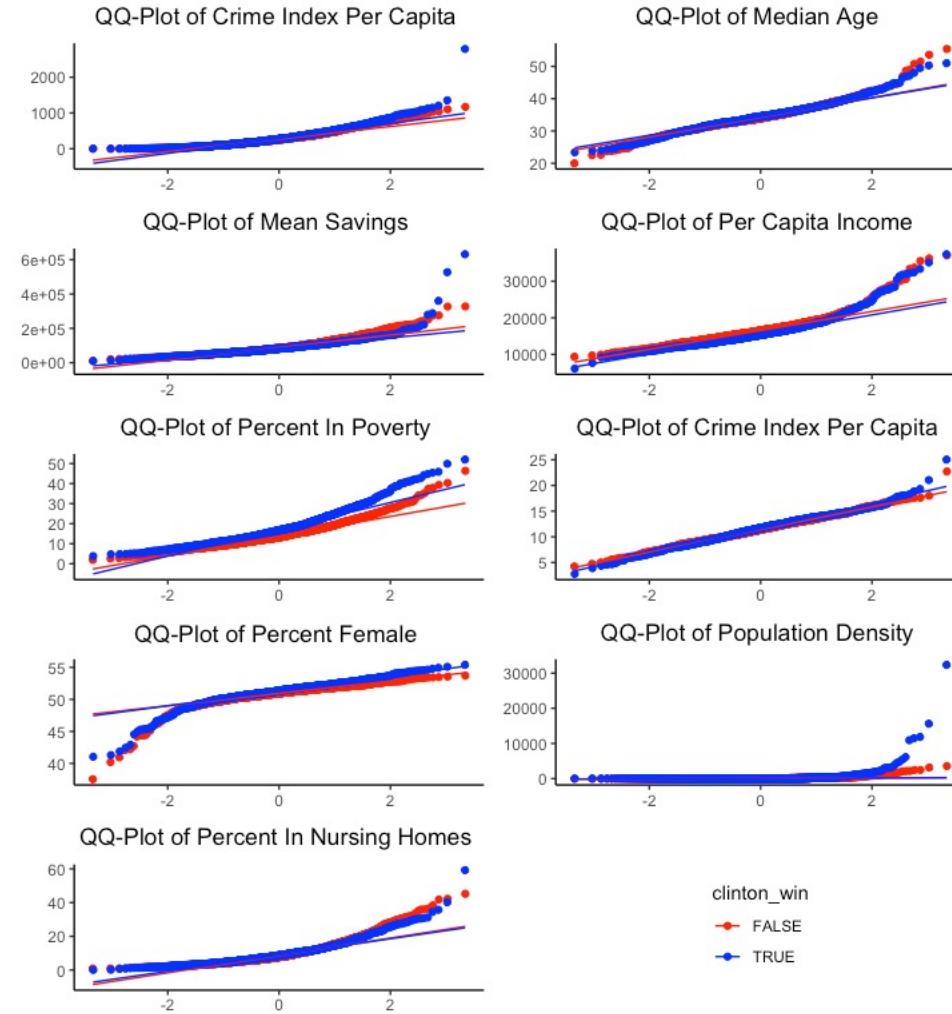


Conclusion

- ▶ There are differences in census measurement that voted for democratic vs. republican
 - ▶ Modeling elections is difficult, but possible
 - ▶ Current baseline: random forest
 - ▶ Future: Time-series model, RNN
 - ▶ Population density, % female, PC income, % in poverty - similar to individual ANOVA results
-
- ▶ Future Work:
 - ▶ Add interaction effects between important variables
 - ▶ Consider statistical models
 - ▶ Time-series models
 - ▶ Ignored cyclical and temporal structure
 - ▶ Optimize Parameters (grid search)
 - ▶ Use my PCA from my report as input to models
 - ▶ Use cross-validation

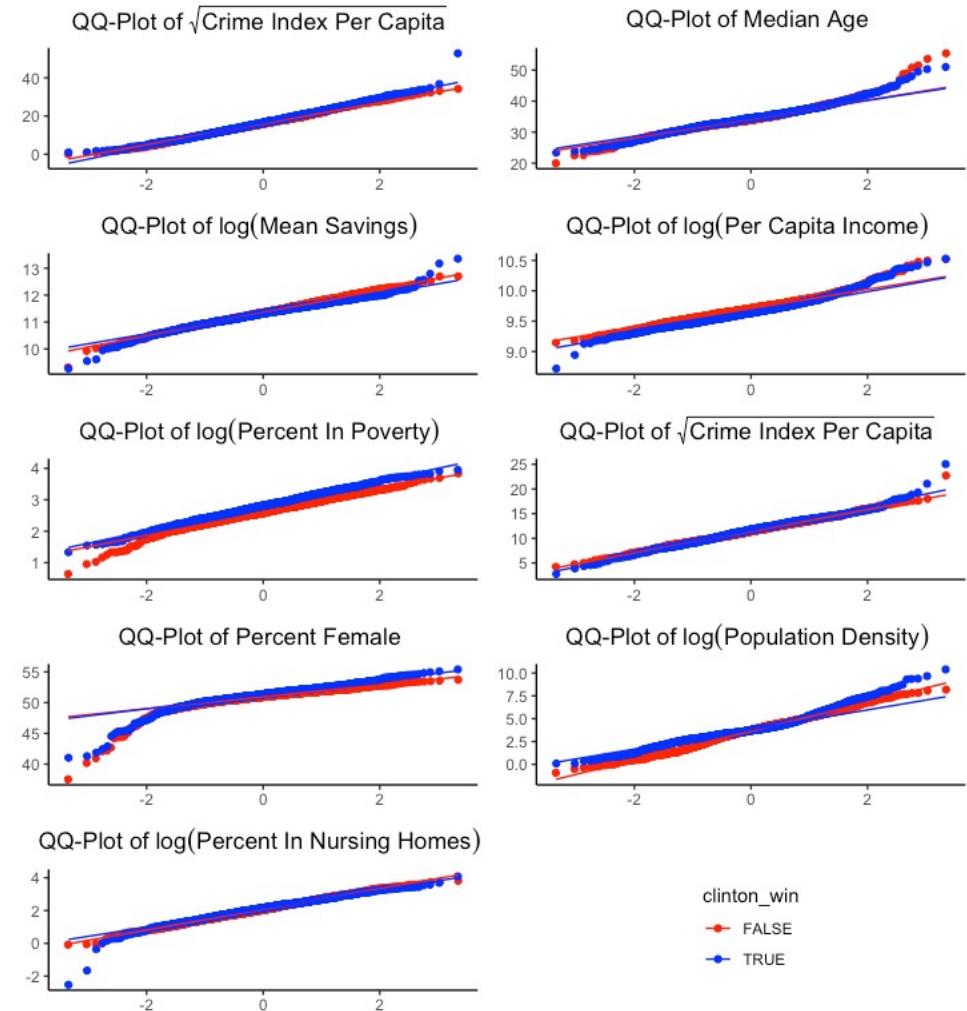
Additional Graphs

► Grouped raw QQ - plots



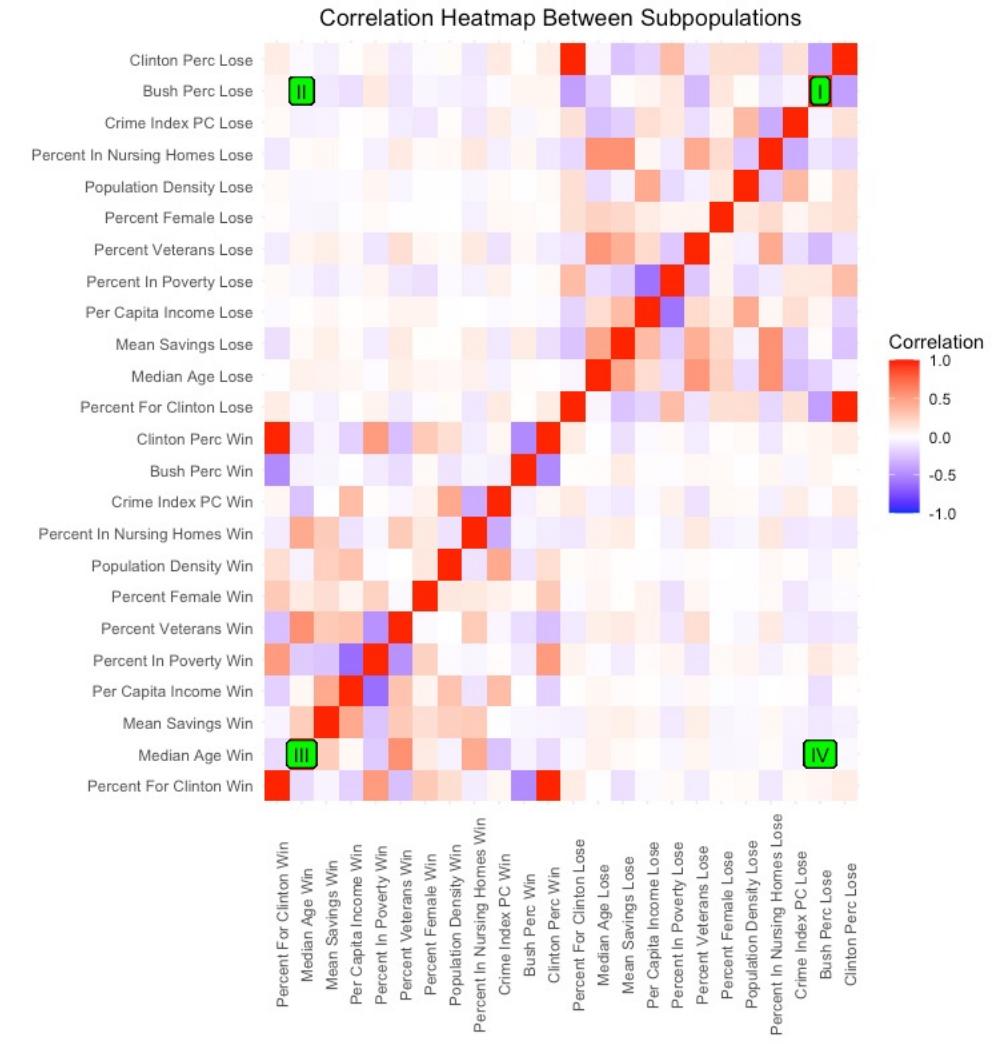
Additional Graphs

- ▶ Grouped transformed QQ - plots



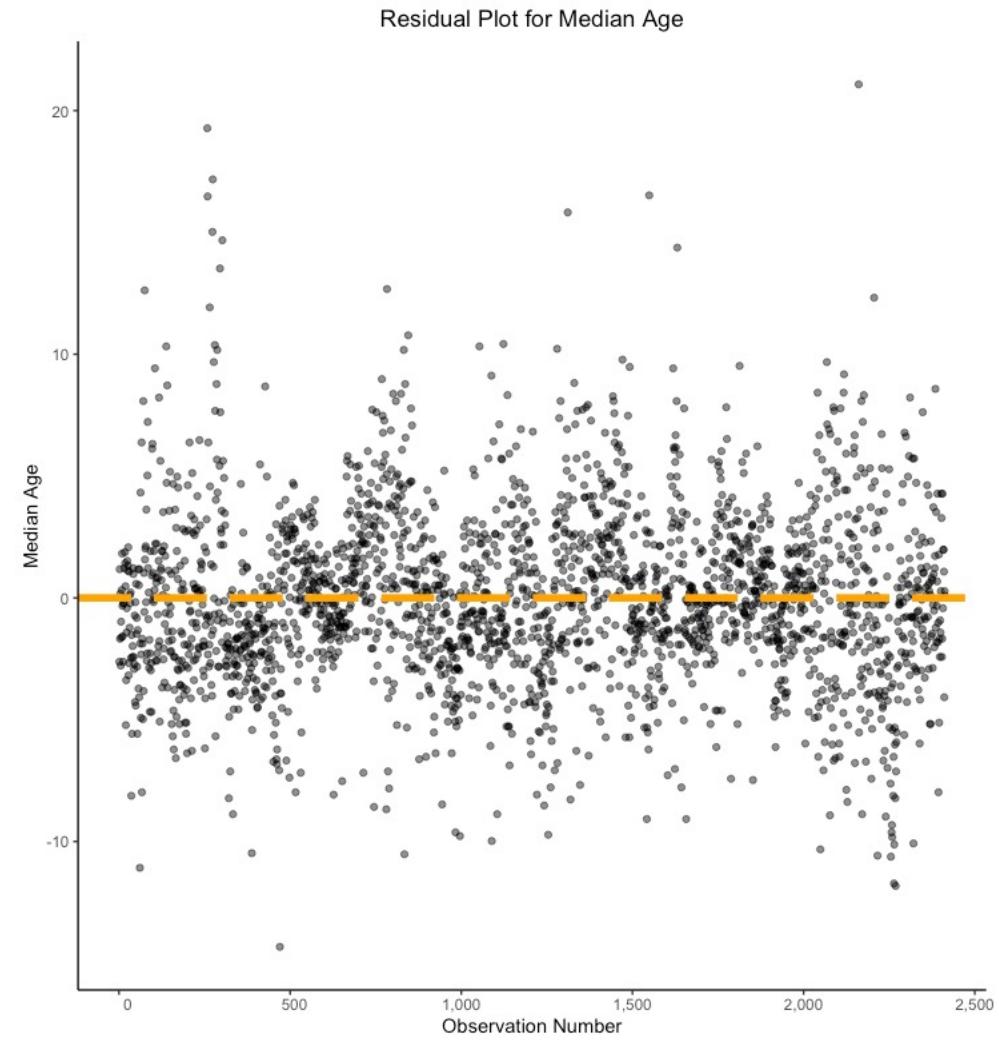
Additional Graphs

► Subpopulation heatmap



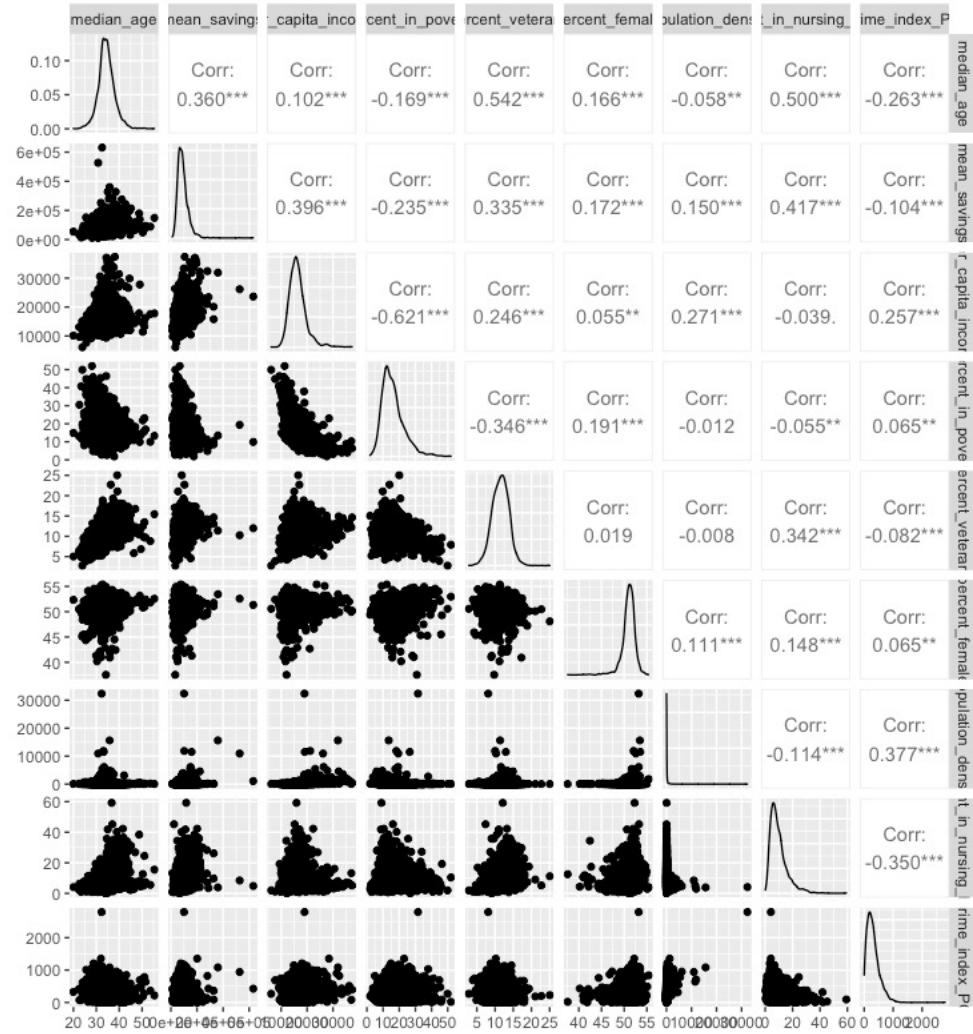
Additional Graphs

- Residual Plot - MANOVA



Additional Graphs

▶ Scatterplot matrix



Confusion Matrix - Random Forest

- ▶ Equally good/bad at classifying/misclassifying election results

```
> conf.matrix  
      predicted  
      observed FALSE TRUE  
      FALSE     177   67  
      TRUE      71  167
```

My R Package VIP List

- ▶ Rvest
- ▶ Dplyr
- ▶ Stringr
- ▶ ggplot2
- ▶ Keras

