

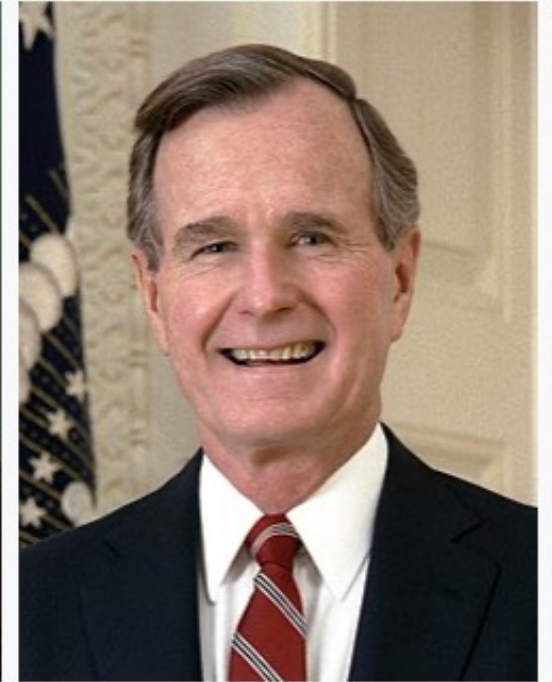
An Exploratory Analysis on US Election Results for 1992 in Preparation for Prediction Tasks

Jared Thacker

Objective

1. Is there a difference in demographic measurements 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. Answer to Question 1
5. Answer to Question 2
6. Answer to Question 3
7. Conclusion
8. Extra Graphs - MANOVA Assumptions
9. Useful Python Packages



Focus of this presentation

Data-Sourcing

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

WIKIPEDIA

The Free Encyclopedia

ArticleTalk

ReadEditView historySearch Wikipedia

</

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

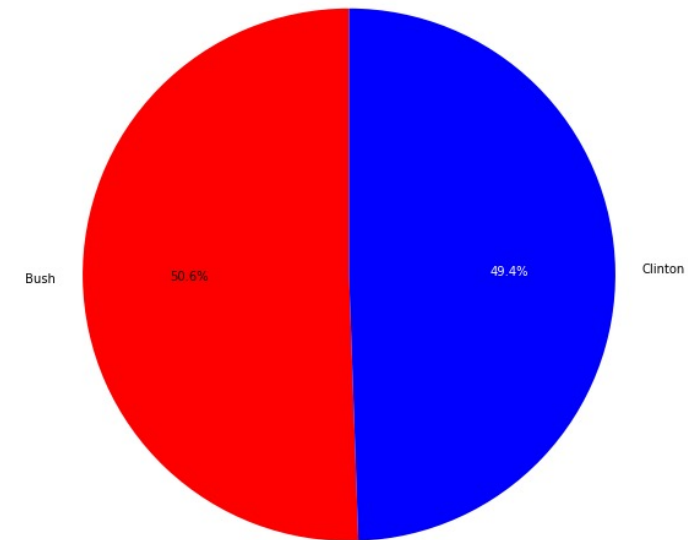
- Clinton Win - Binary, nominal
- 2413 Observations (counties)
- 9 predictors, 1 response
- 1220 counties that Bush won
- 1193 counties that Clinton won
- Subpopulations almost equivalent

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

Response Variable - Class Imbalance?

- ▶ 1220 counties that Bush won
- ▶ 1193 counties that Clinton won

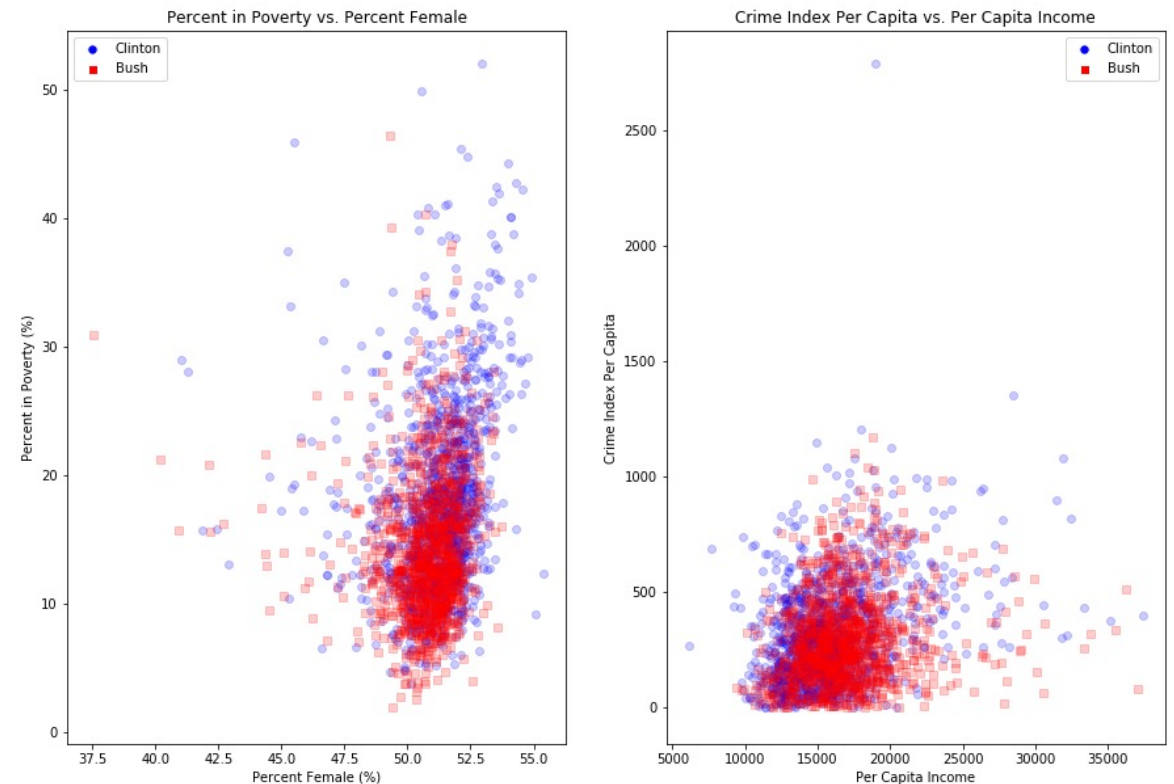
Response Class Breakdown by Clinton Victory or Loss



Exploratory Data Analysis (EDA)

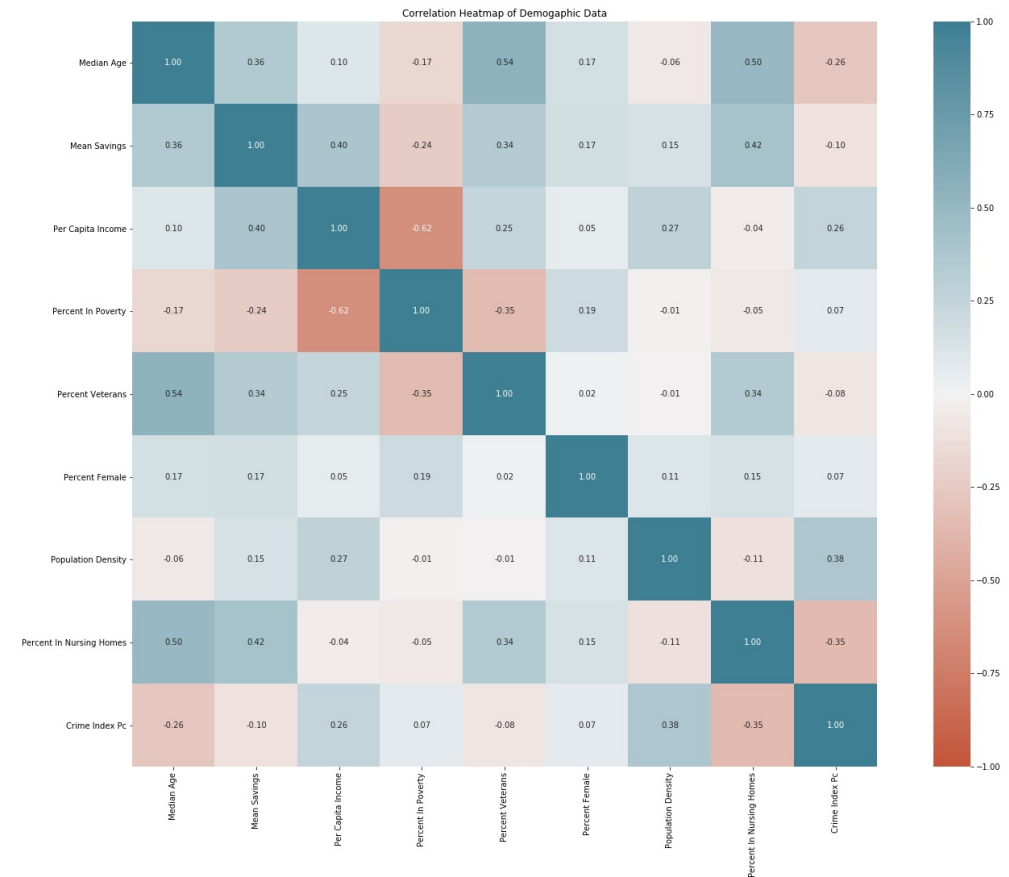
- ▶ From the scatterplots - no clear decision boundary
- ▶ Nonlinear classifiers will be needed

Class separation:
NONLINEAR



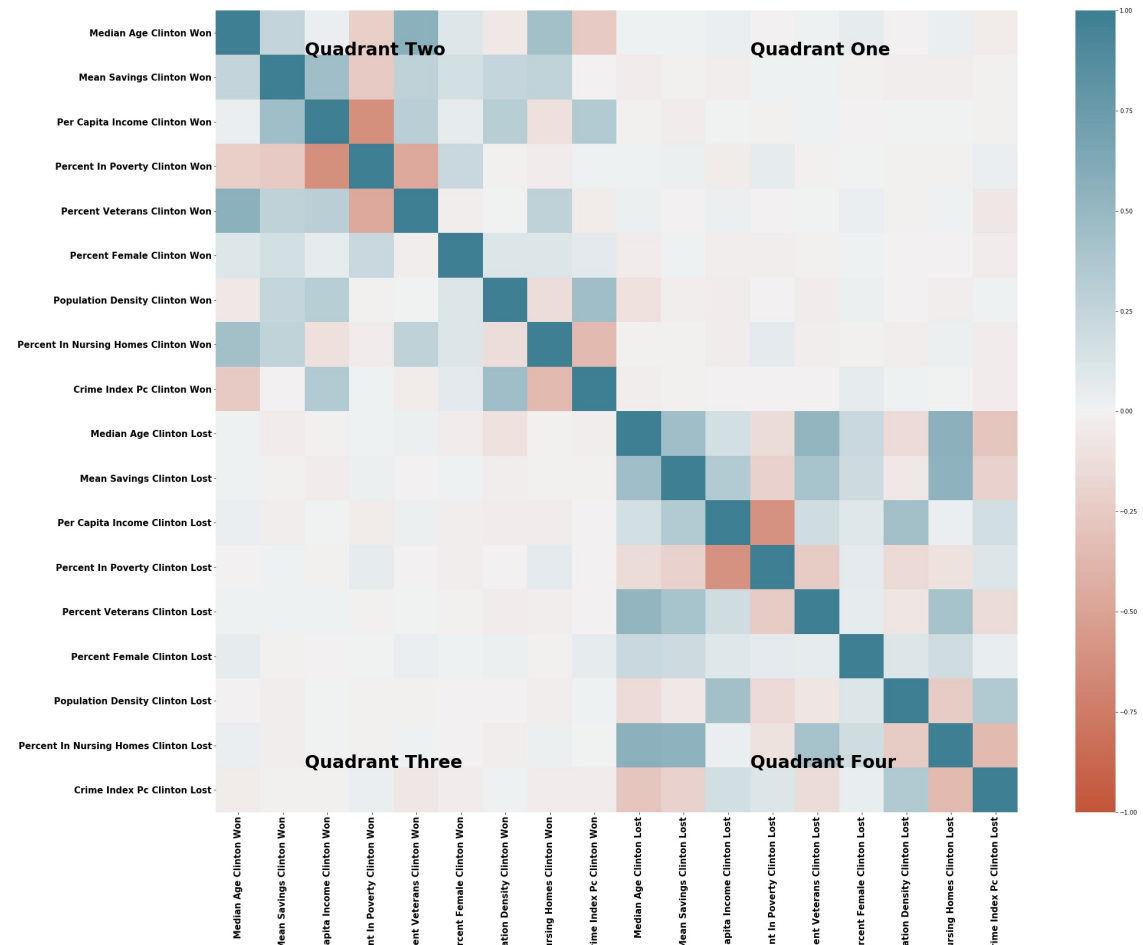
More EDA

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



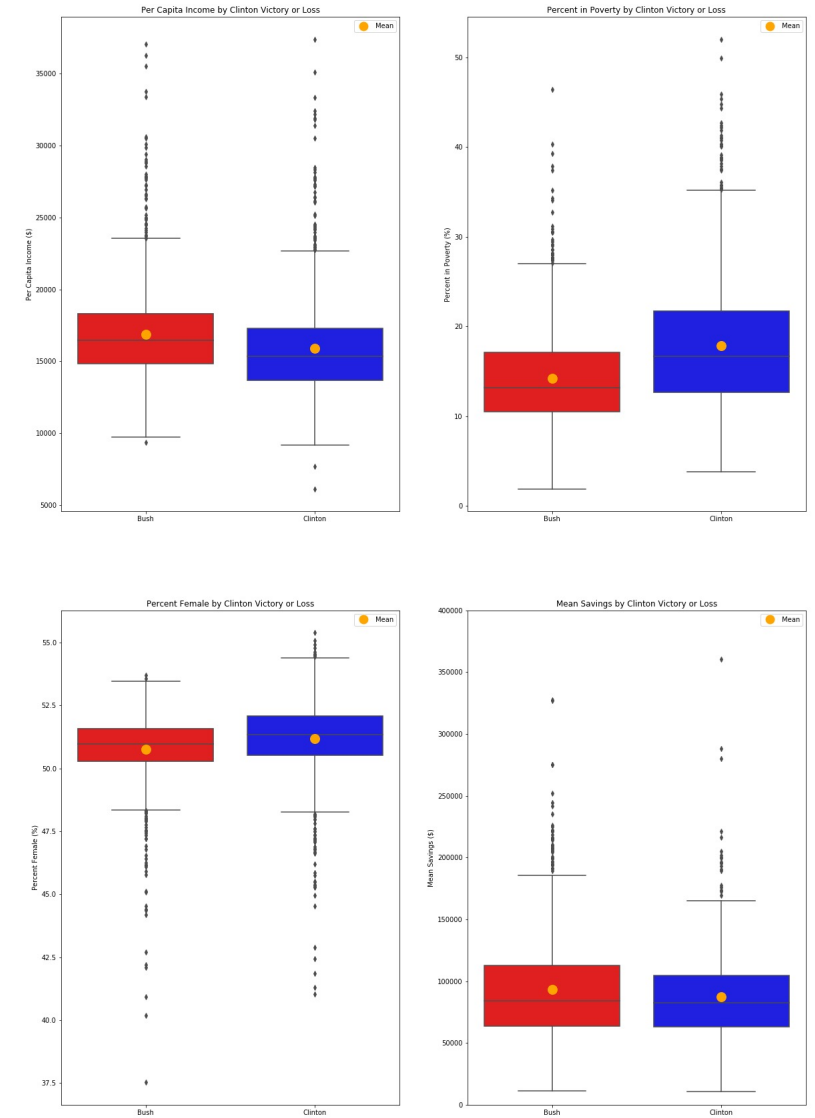
EDA - Between Sample Dependence

- ▶ MANOVA assumption - samples from different populations are independent
- ▶ The assumption holds for linear correlation holds
- ▶ If there is nonlinear correlation, then it will show up in MANOVA residuals and we can trouble shoot from there



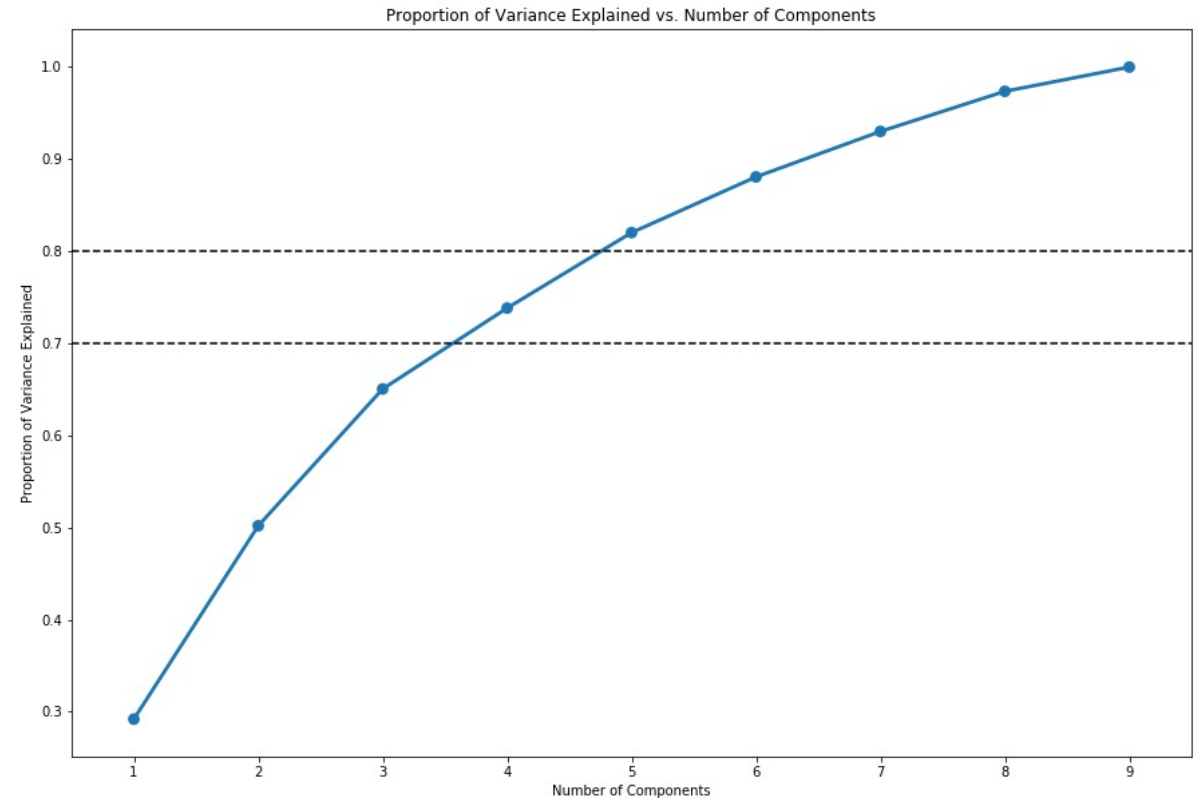
More EDA

- Appearance: small difference between different census measure
- Large sample size -> more power
- If assumptions of MANOVA hold, then the null hypothesis will likely be rejected due to large sample size



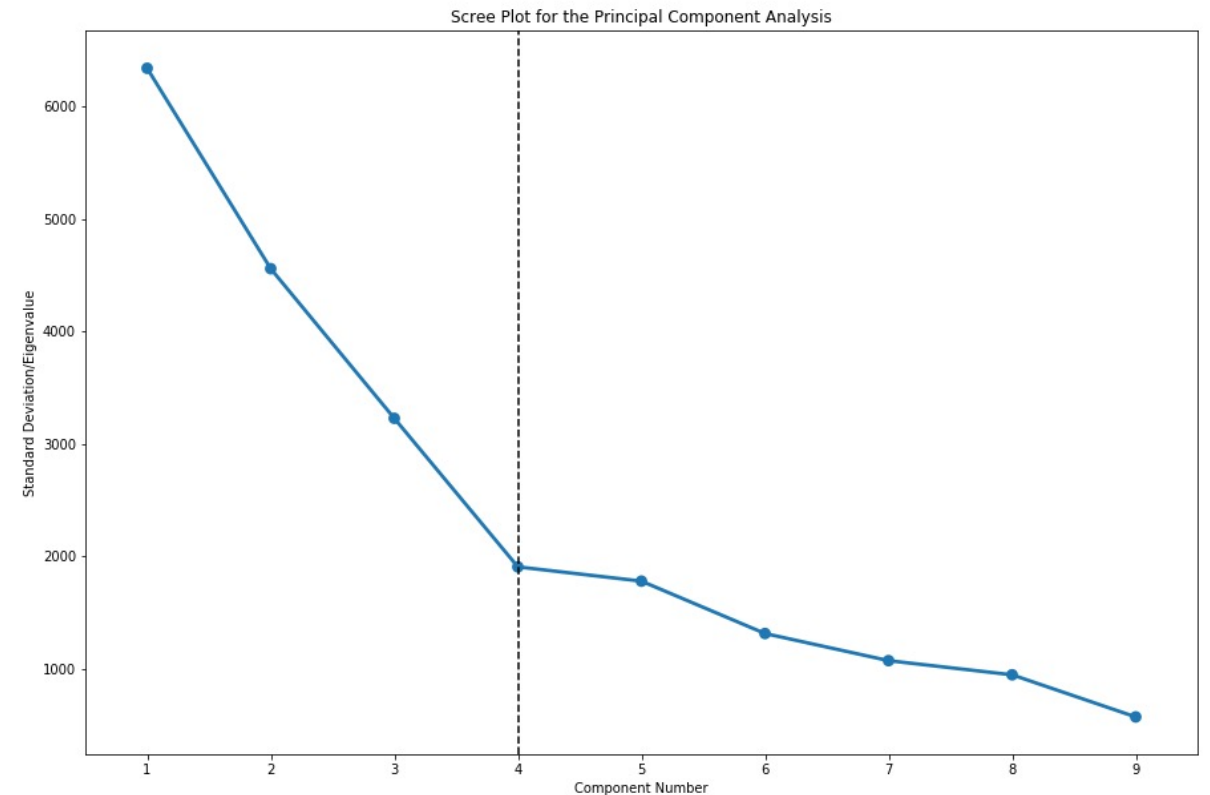
PCA - How Many Components

- ▶ 4th component - 75% of the variance explained
- ▶ Since 70% > of variance explained by 4 components, retain only 4 components



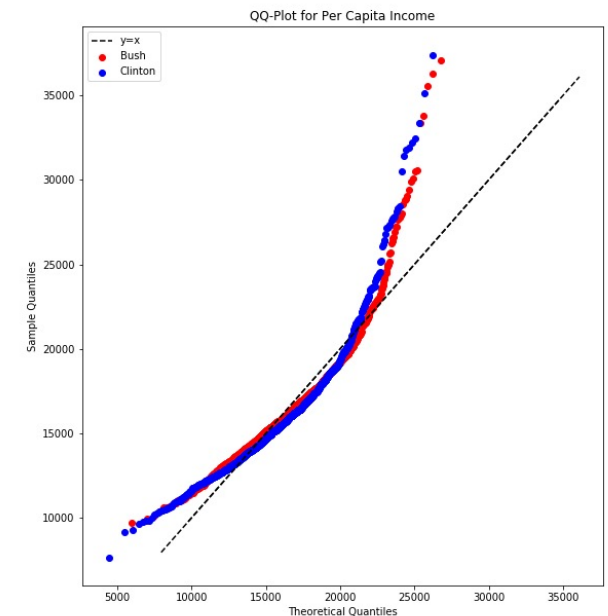
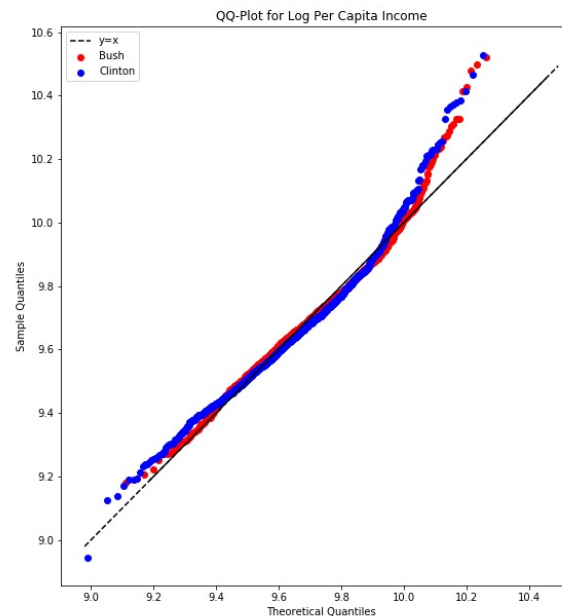
PCA - How Many Components

- ▶ Scree plot - Find the elbow in the graph and retain that number of components
- ▶ Elbow happens at 4 components - retain 4 components



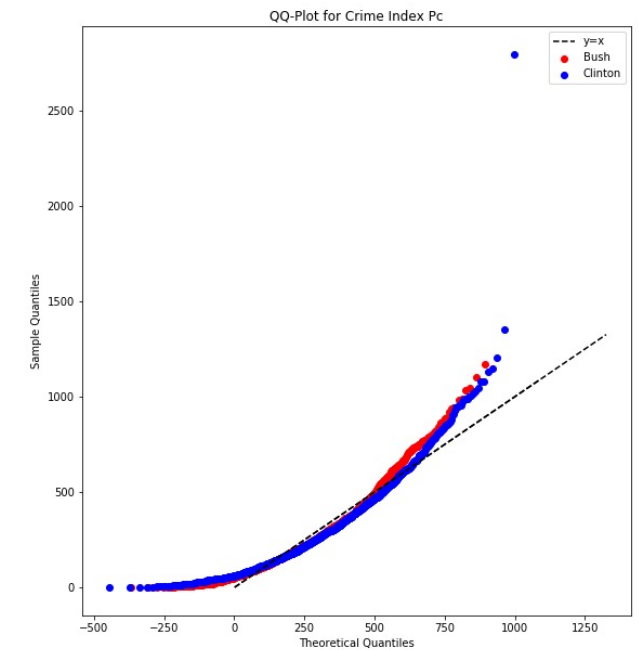
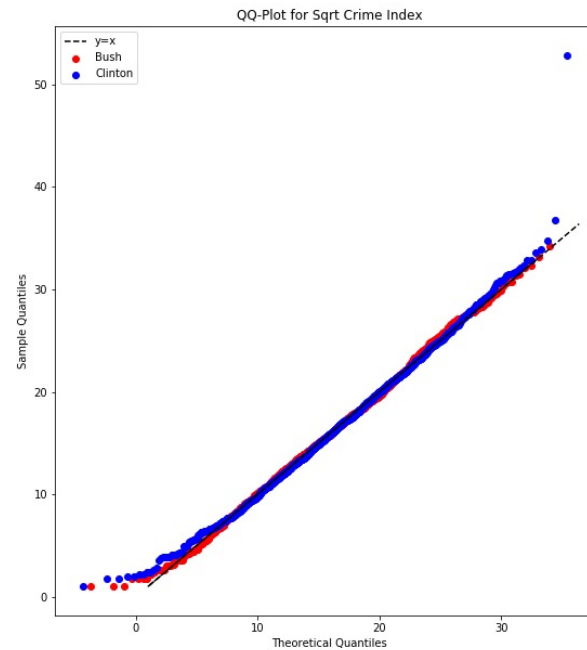
Distribution of the Subpopulations

- ▶ Untransformed variables -> most of the non-normal with 1-2 exceptions
- ▶ Transformed variables -> approximately normal or at least do not stray far from normality
- ▶ Per Capita Income - Not normal
- ▶ Log(Per Capita Income) - Normal



Distribution of the Subpopulation (continued)

- ▶ Crime Index Per Capita - not normal
- ▶ Sqrt(Crime Index PC) - approximately normal



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