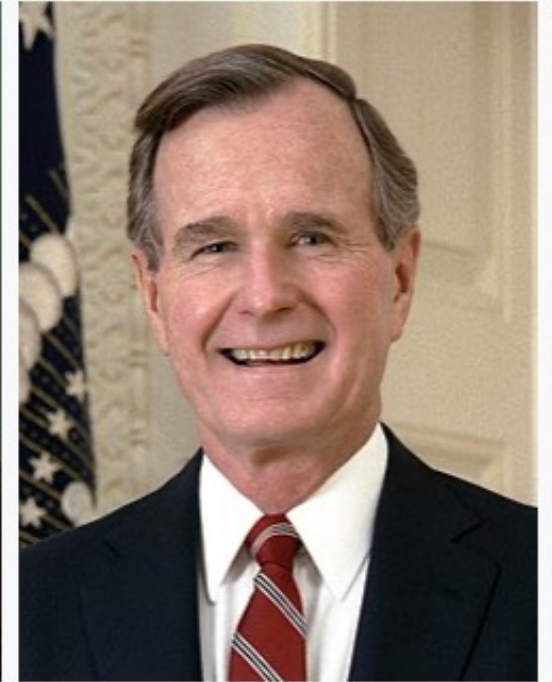# A Multivariate Analysis on U.S Presidential Election Results From 1992

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 ignoring the temporal structure? Baseline model?

3. Which variables are the most important?

Democratic Candidate: Bill Clinton

Republican Candidate: George H. W. Bush

# Data-Sourcing

- Two Sources
  - Ufl.edu – census variables by county
  - Wikipedia – election results by county



An example of a Wikipedia table that was scraped

# Exploratory Data Analysis (EDA)



Percent Female vs. Percent in Poverty

Class separation: NONLINEAR



Vote Percentage for Clinton by State

# More EDA

▶ Moderately strong correlation between variables

▶ PCA might be appropriate
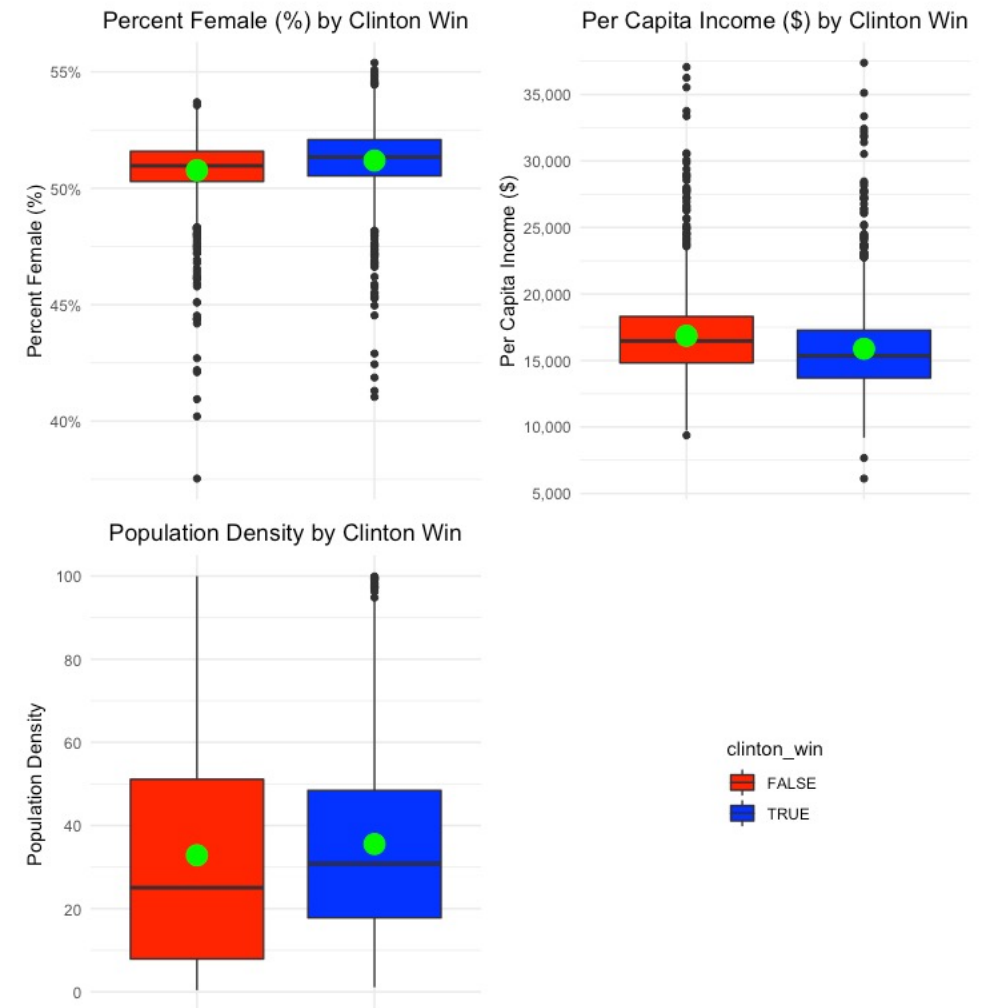


Correlation Heatmap for the Whole Dataset

# 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: $p < 0.0001$

# 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%
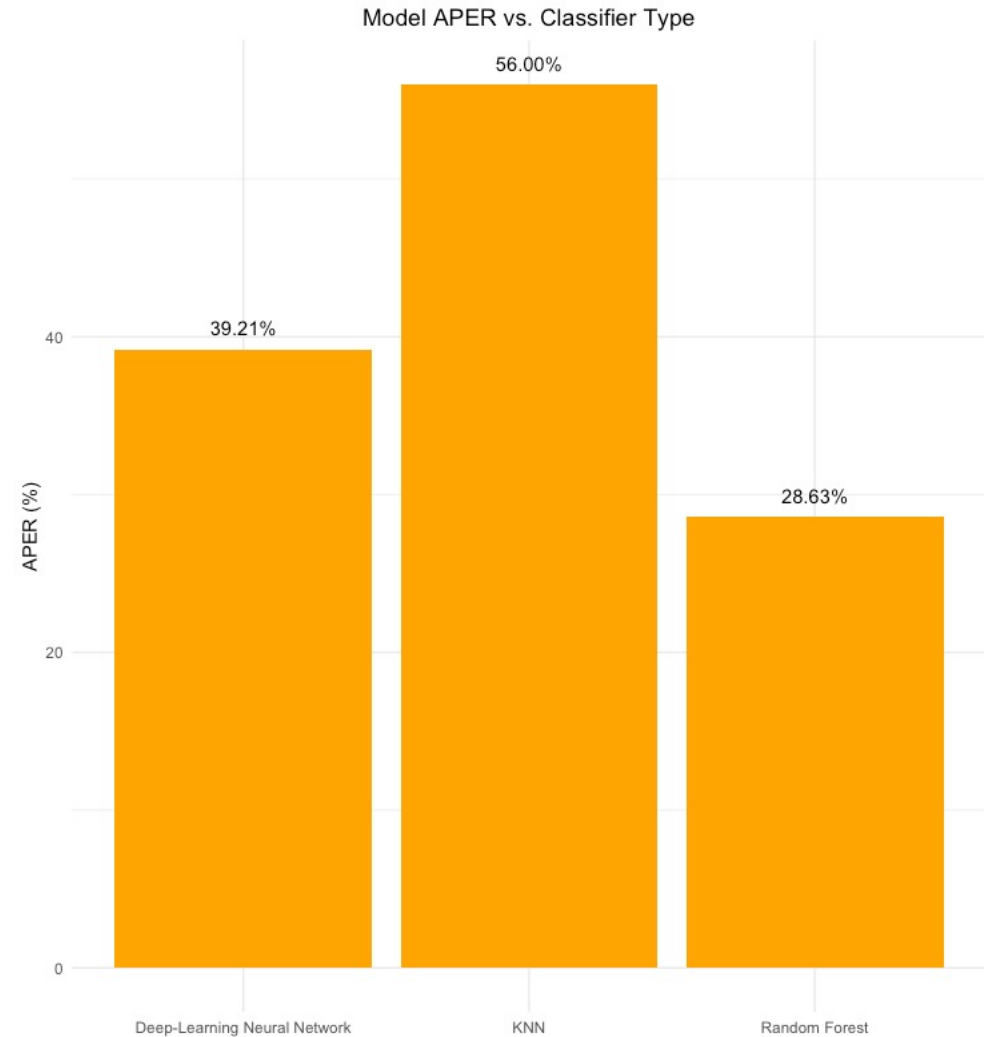  - ▶ Activation function: Sigmoid
  - ▶ Loss function = "binary cross entropy"
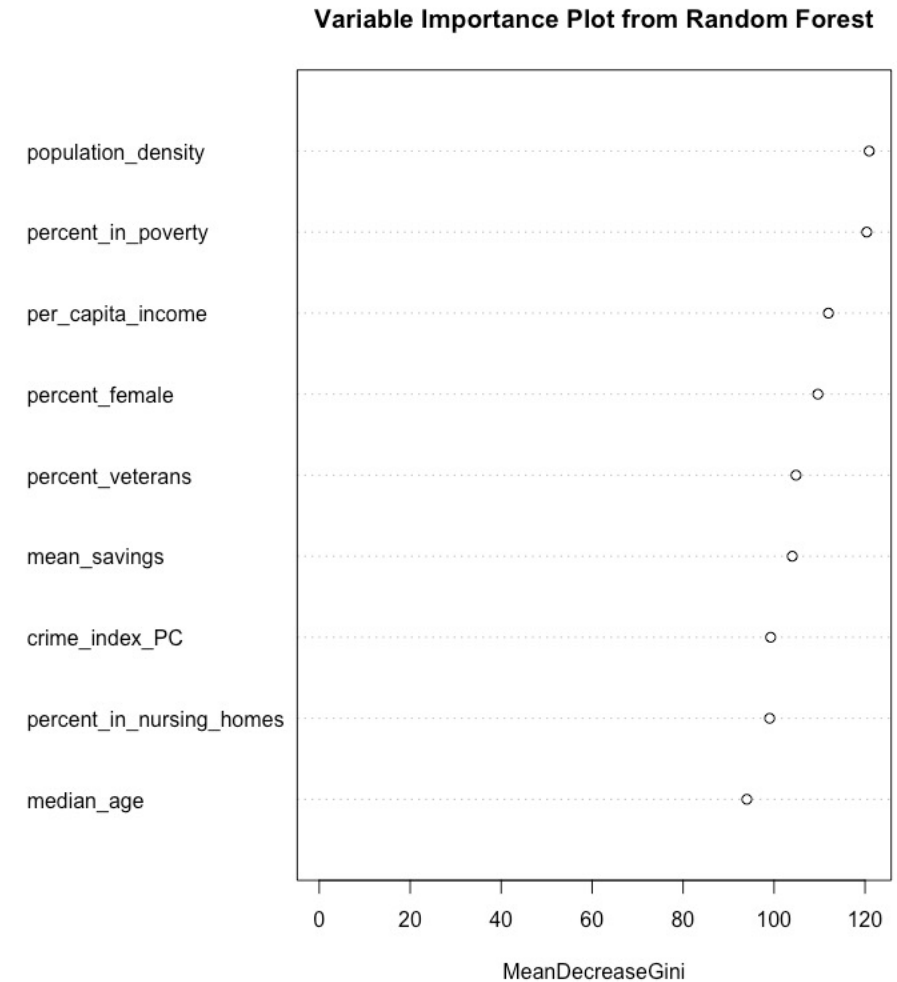  - ▶ Optimizer = "RMSprop"

# Question 2: Predictive Modeling

- ▶ Random forest – highest performer
- ▶ Neural Network
  - ▶ Not enough data (inconclusive)
- ▶ There are more ML models
  - ▶ I chose just three
- ▶ I should've considered statistical models
  - ▶ ML offers no advantage sometimes
- ▶ We're ignoring the temporal and cyclical nature of election cycles – this is the future of this study (Time series models, RNNs)

Model APER vs. Classifier Type

# Question 3: Variable Importance

▶ Random forest – highest performer: Population Density

▶ Random forest importance is *the same as* ANOVA results
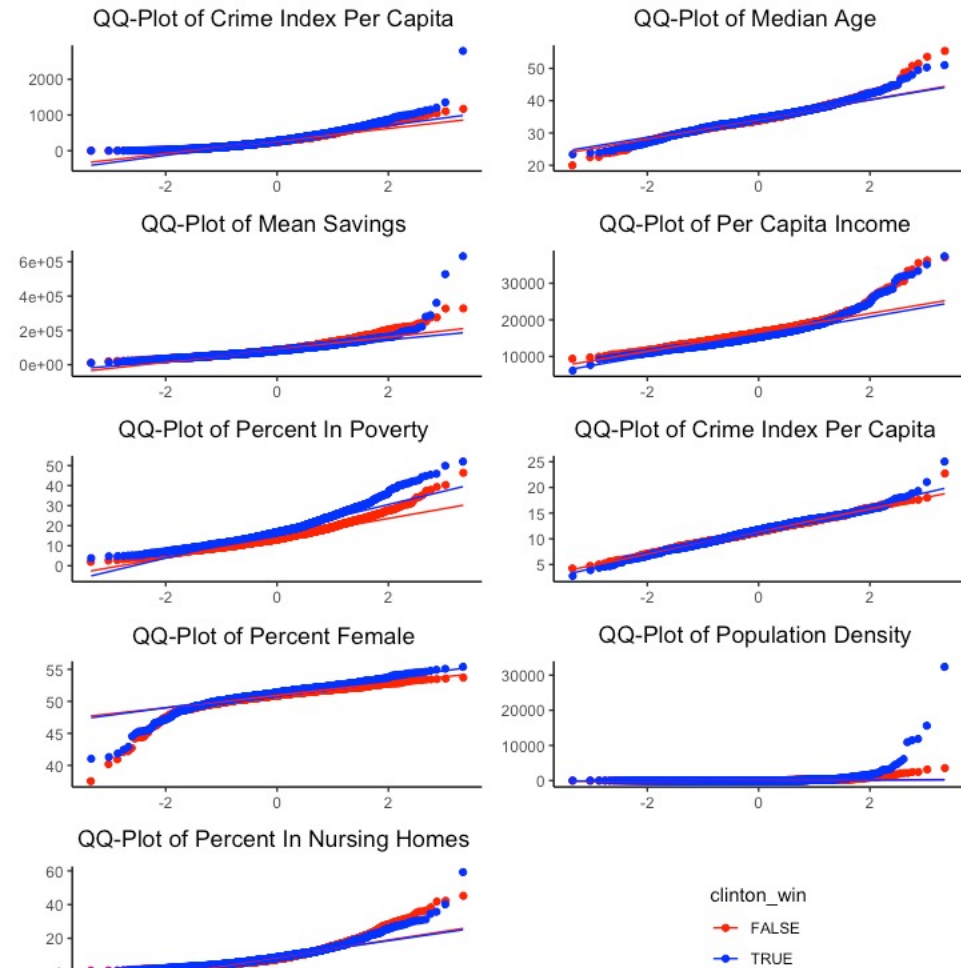
**Variable Importance Plot from Random Forest**

# 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 – same as individual ANOVA results

- Future Work:
  - Add interaction effects between important variables
  - Time-series models
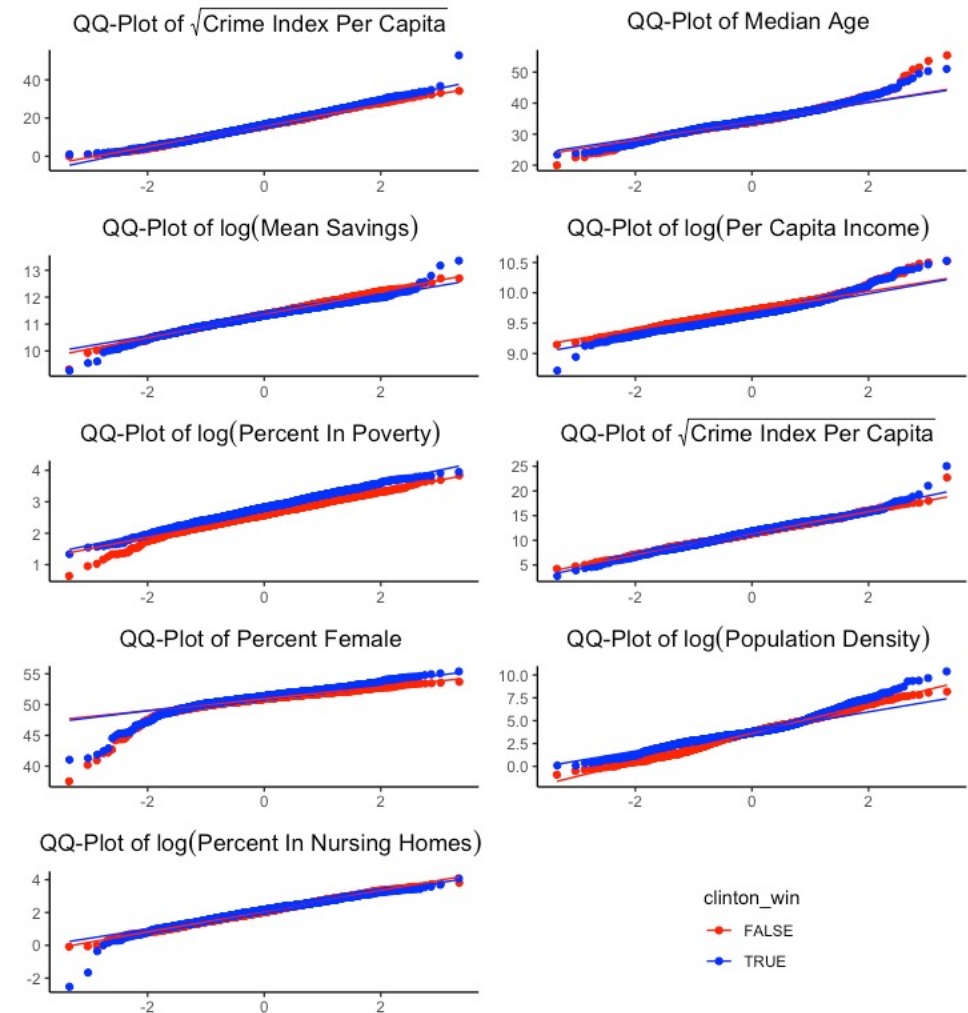  - Optimize Parameters (grid search)
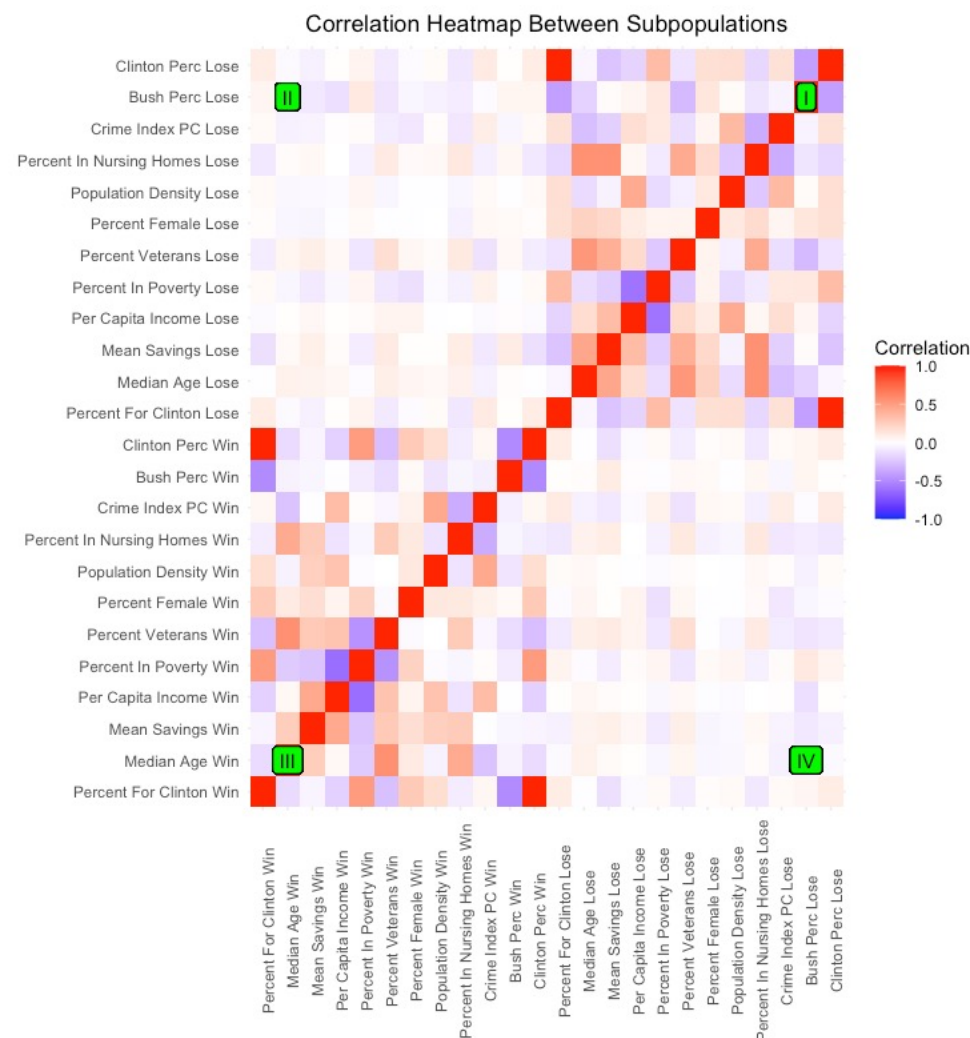  - Use cross-validation

# Additional Graphs

▶ Grouped raw QQ - plots

# Additional Graphs
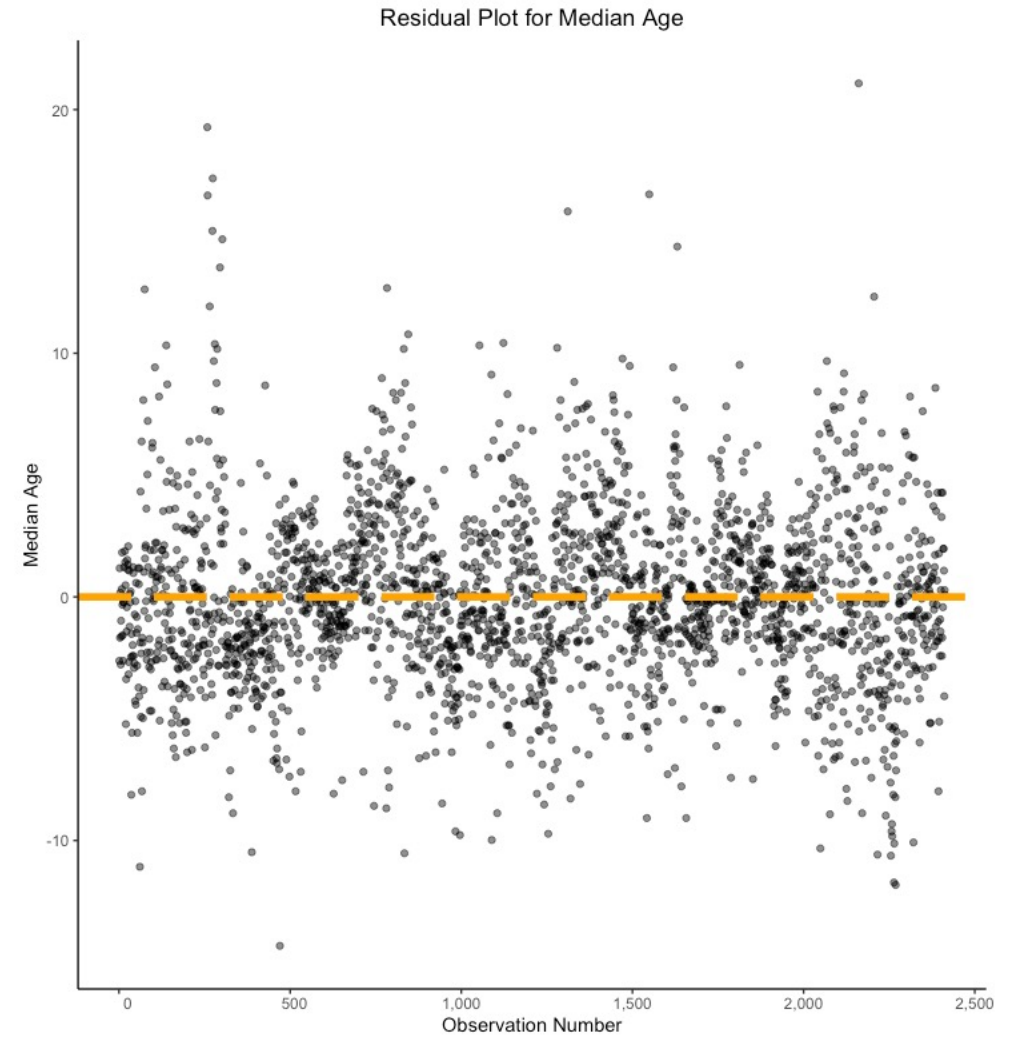
▶ Grouped transformed QQ - plots
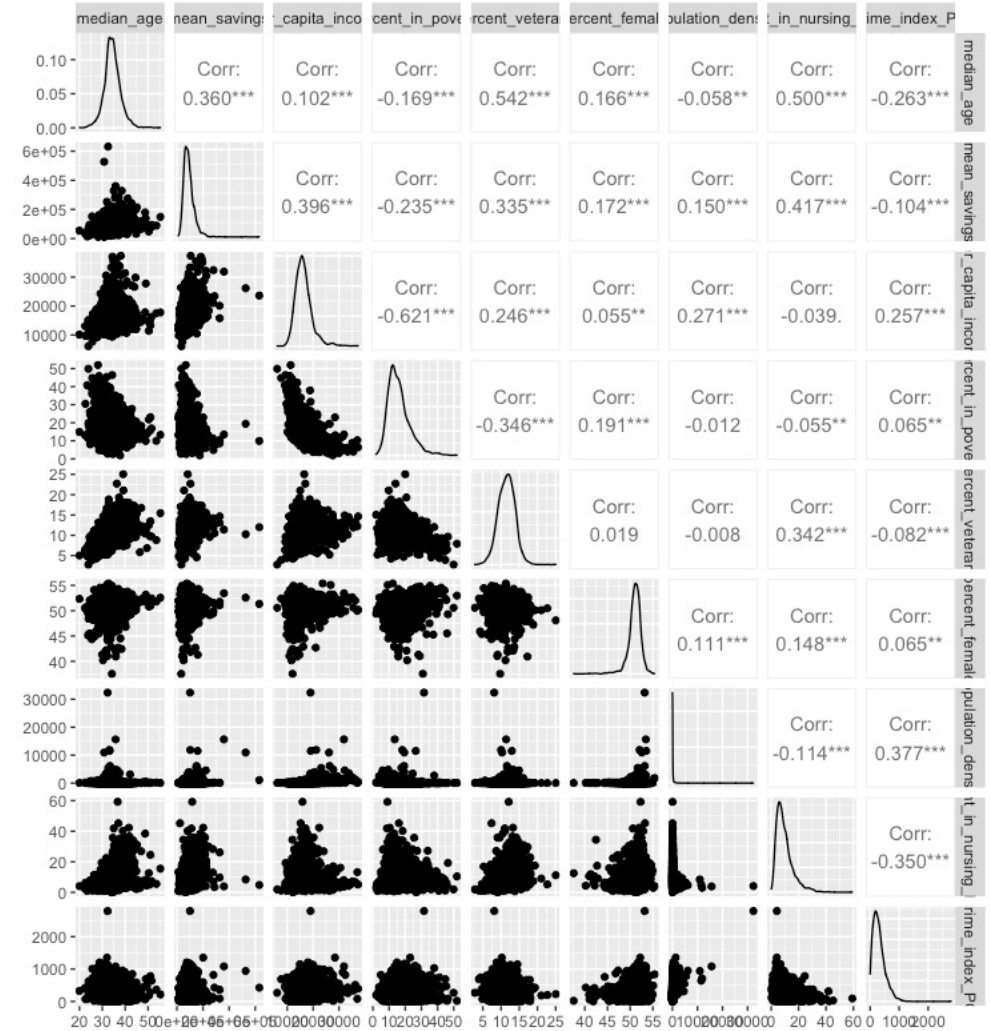
# Additional Graphs

▶ Subpopulation heatmap

# Additional Graphs

▶ Residual Plot - MANOVA

# Additional Graphs

▶ Scatterplot matrix

# My R Package VIP List

▶ Rvest

▶ Dplyr

▶ Stringr

▶ ggplot2

▶ Keras