# Household Conditions by Geographic School District

Jon Geiger, Noel Goodwin, Abigail Joppa

Week 7/8

# Research Question

We know that correlations between graduation rate and household/demographic conditions vary by region.

**What household conditions are the biggest indicators for graduation rates across school districts?**

How does this differ across:

➢ Regions?

➢ Tiers of school district funding?

**Does district-wide assessment data provide a significant improvement to a regression model?**

# Last Week's Steps..

Visualize predictor distributions for regression

Improve regression model predicting graduation rates using HH conditions and race

Continue to look for prior research that looks at graduation rates and indicators
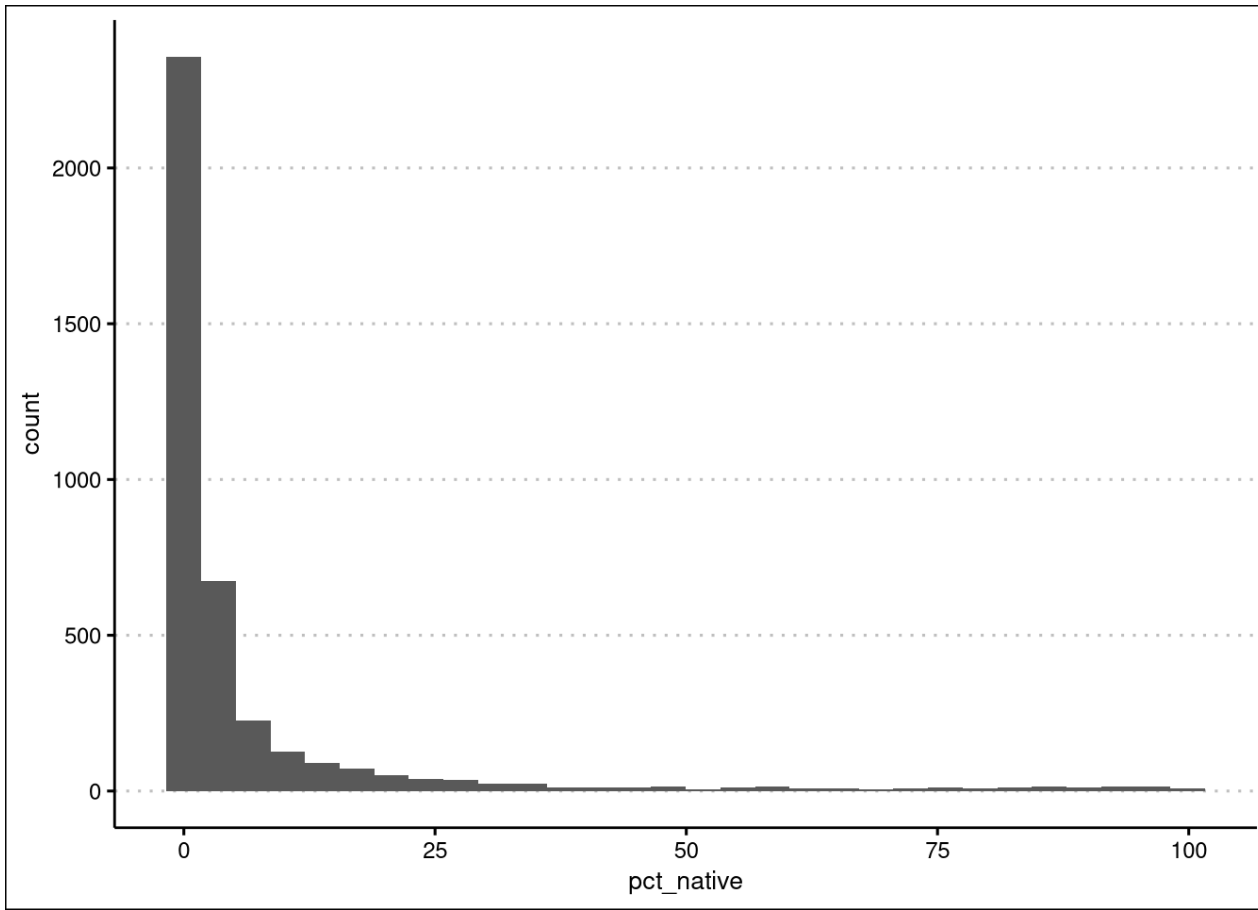
# Last Week's Difficulties:

- Difficult to "group by" racial/ethnic groups

- Defining regions

- Household income group data
  - Proxy with school district funding?
  - Look not at income groups, but ranked groups of amount of funding n

# Difficult to "group by" racial/ethnic groups

- Got rid of PI
  - Just over 800 districts with PI students
  - Of those 800 districts, the majority have under 5 percent

- Native racial group

# Native Population



• We did not want to disregard the native population, so we turned our percent native variable into a dummy variable

  • If native students comprise most of the district population, we might assume that the district is on a reservation (proxy)
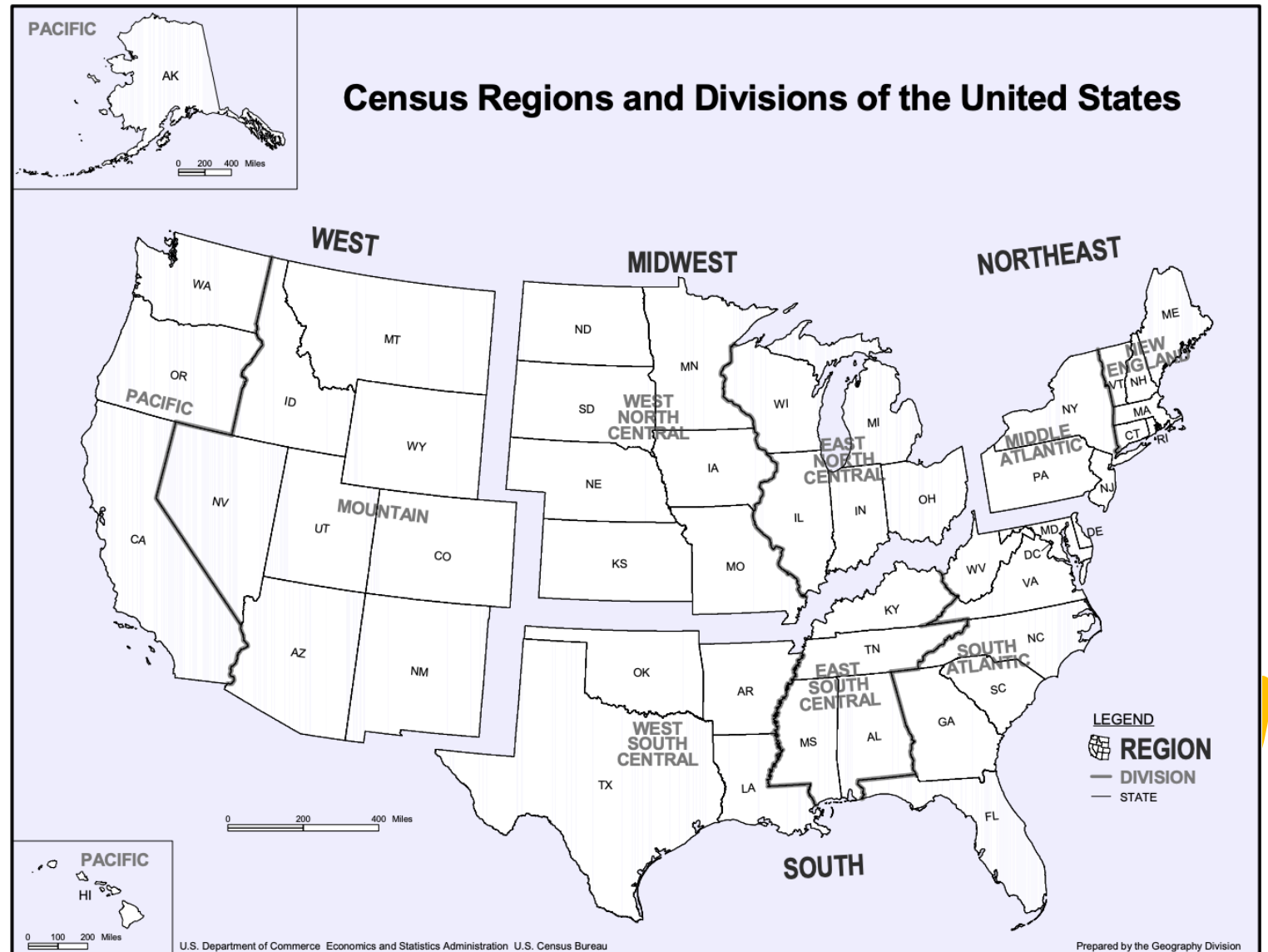
# Difficult to "group by" racial/ethnic groups continued...

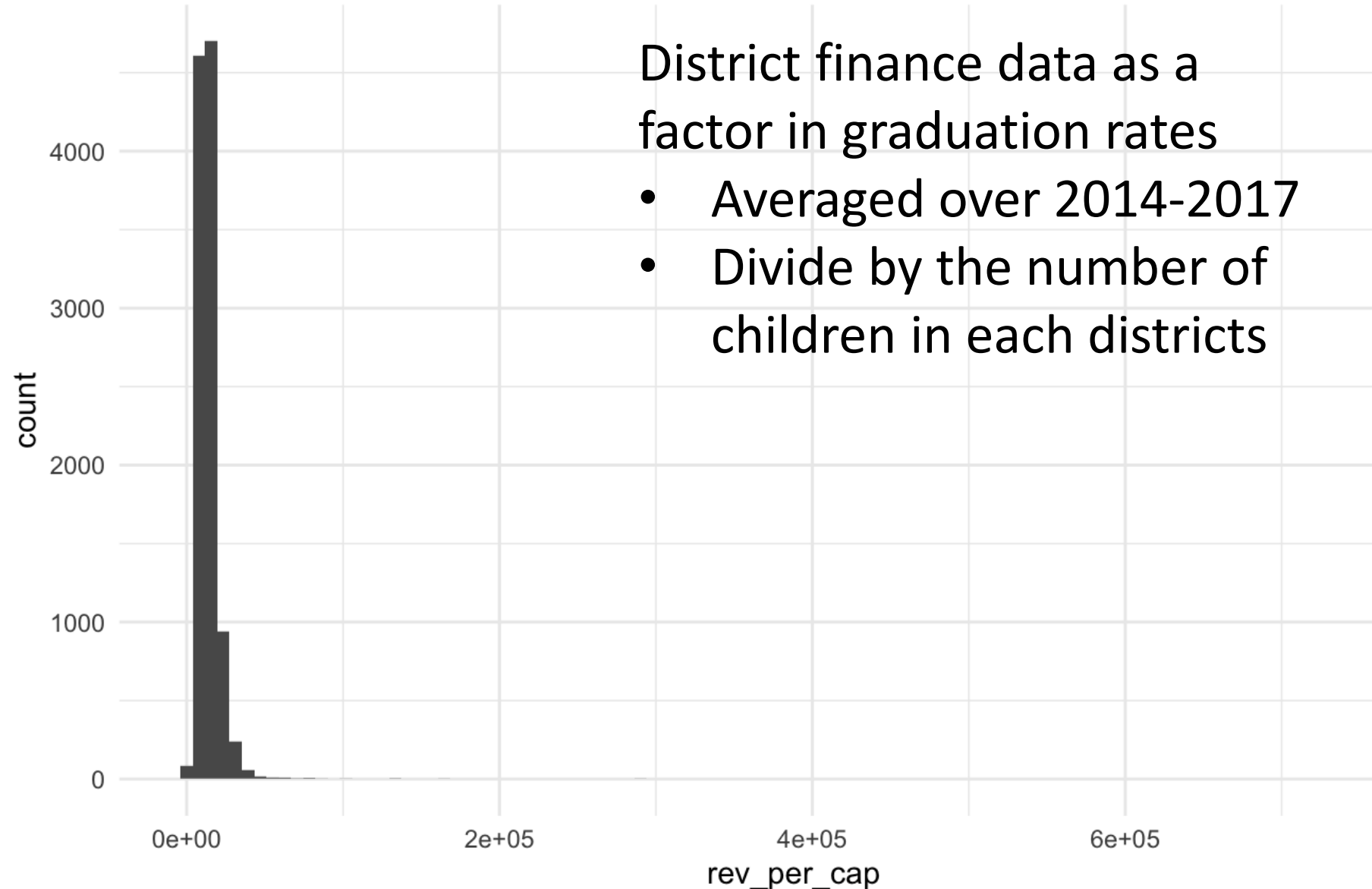- Created a column for each district that indicates the majority racial group

| | predom_race | n |
|---|---|---|
| | <chr> | <int> |
| 1 | Asian | 52 |
| 2 | Black | 531 |
| 3 | Hispanic/Latino | 1213 |
| 4 | Native American | 198 |
| 5 | White | 9916 |
| 6 | NA | 1404 |

| dist | children | predom_race |
|---|---|---|
| Fort Rucker School District | 985 | White |
| Maxwell AFB School District | 292 | White |
| Albertville City School District | 4591 | White |
| Marshall County School District | 8299 | White |
| Hoover City School District | 15397 | White |
| Madison City School District | 9416 | White |
| Leeds City School District | 2324 | White |
| Boaz City School District | 1644 | White |
| Trussville City School District | 4476 | White |
| Alexander City City School District | 3064 | White |
| Andalusia City School District | 1425 | White |
| Arab City School District | 1454 | White |
| Athens City School District | 3824 | White |
| Attalla City School District | 1144 | White |
| Saraland City School District | 2728 | White |
| Satsuma City School District | 922 | White |
| Alabaster City School District | 6469 | White |
| Pelham City School District | 4044 | White |
| Pike Road City School District | 1813 | White |

# Defining Regions



Census Regions and Divisions of the United States

# School District Funding



District finance data as a factor in graduation rates

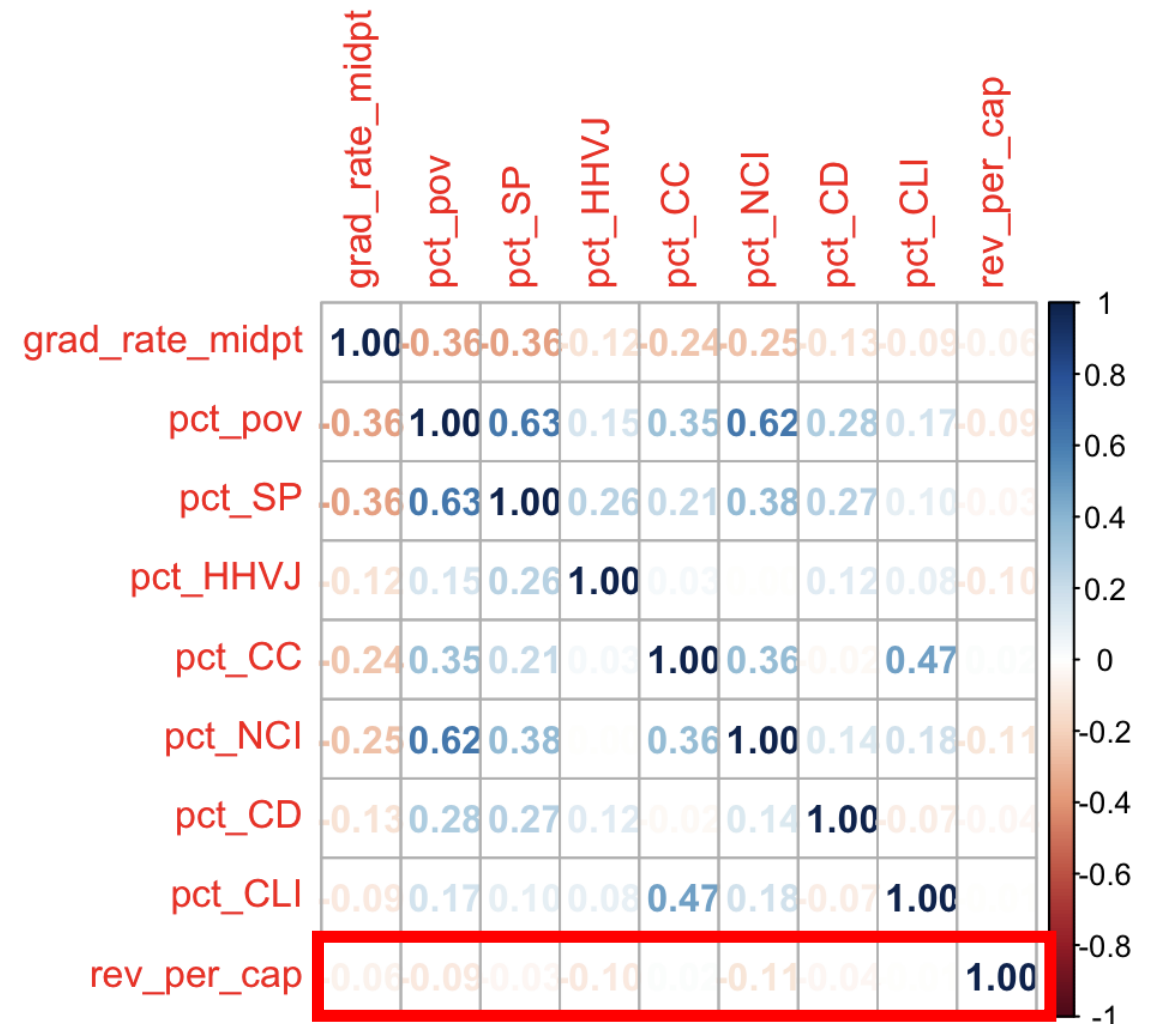- Averaged over 2014-2017
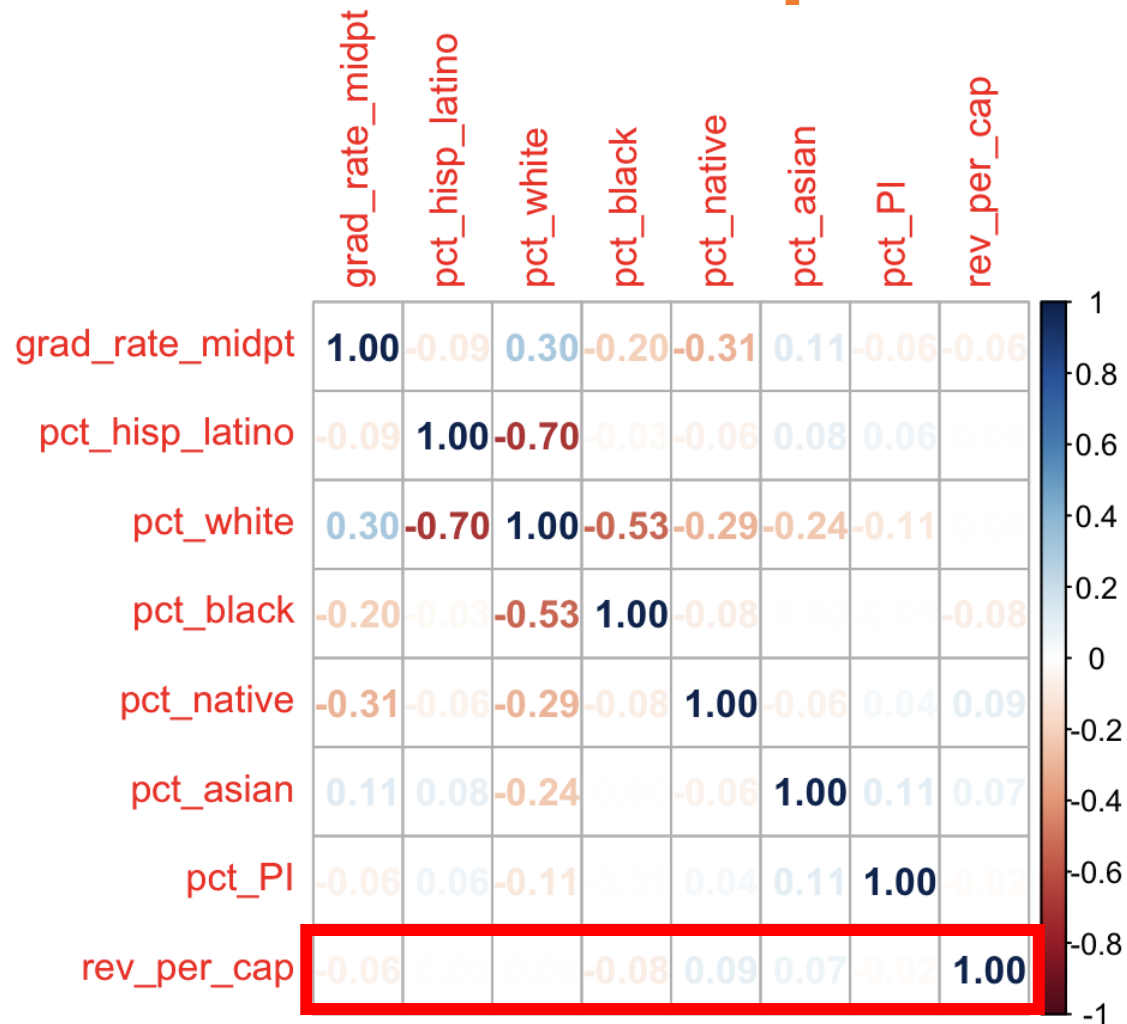- Divide by the number of children in each districts
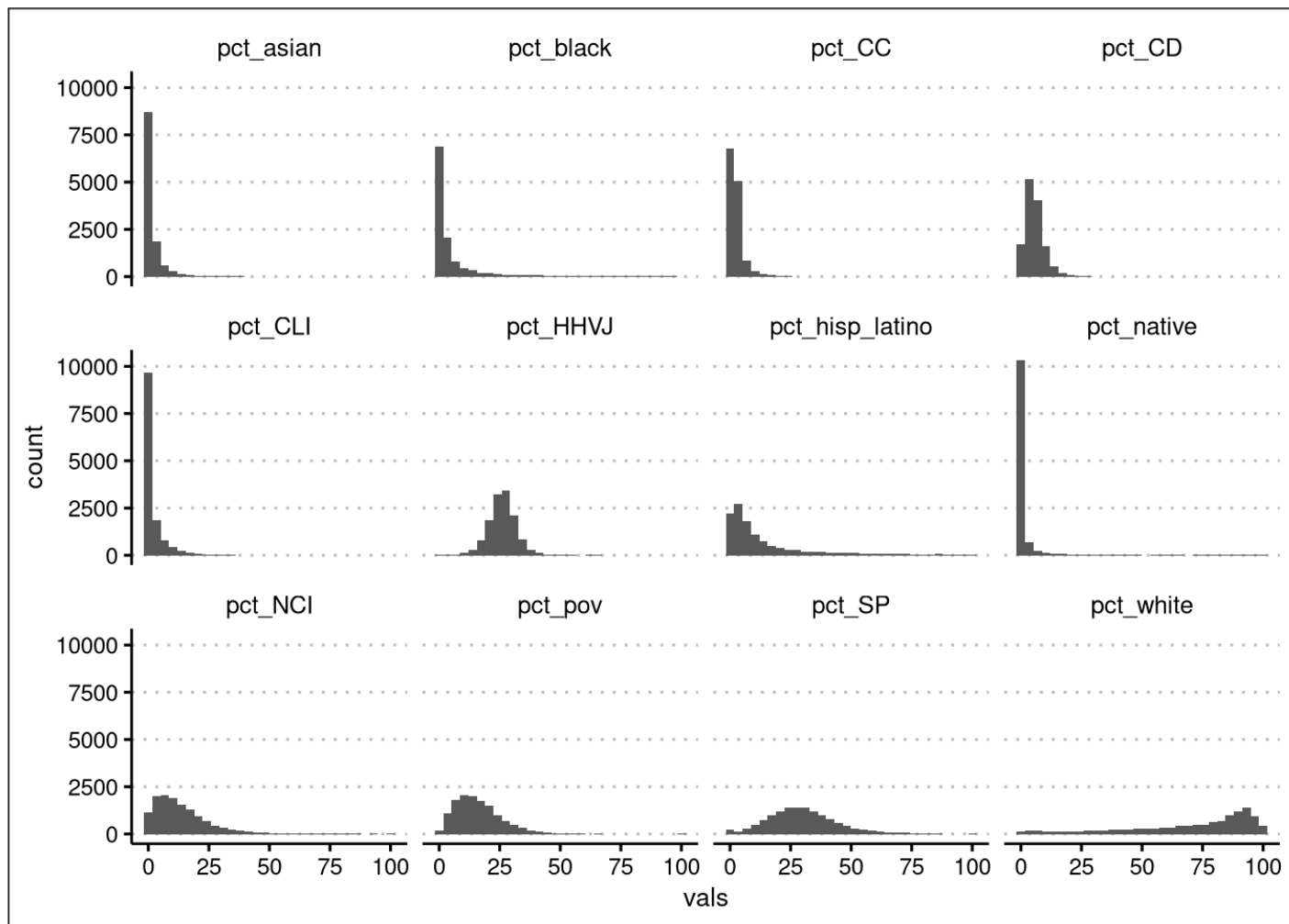
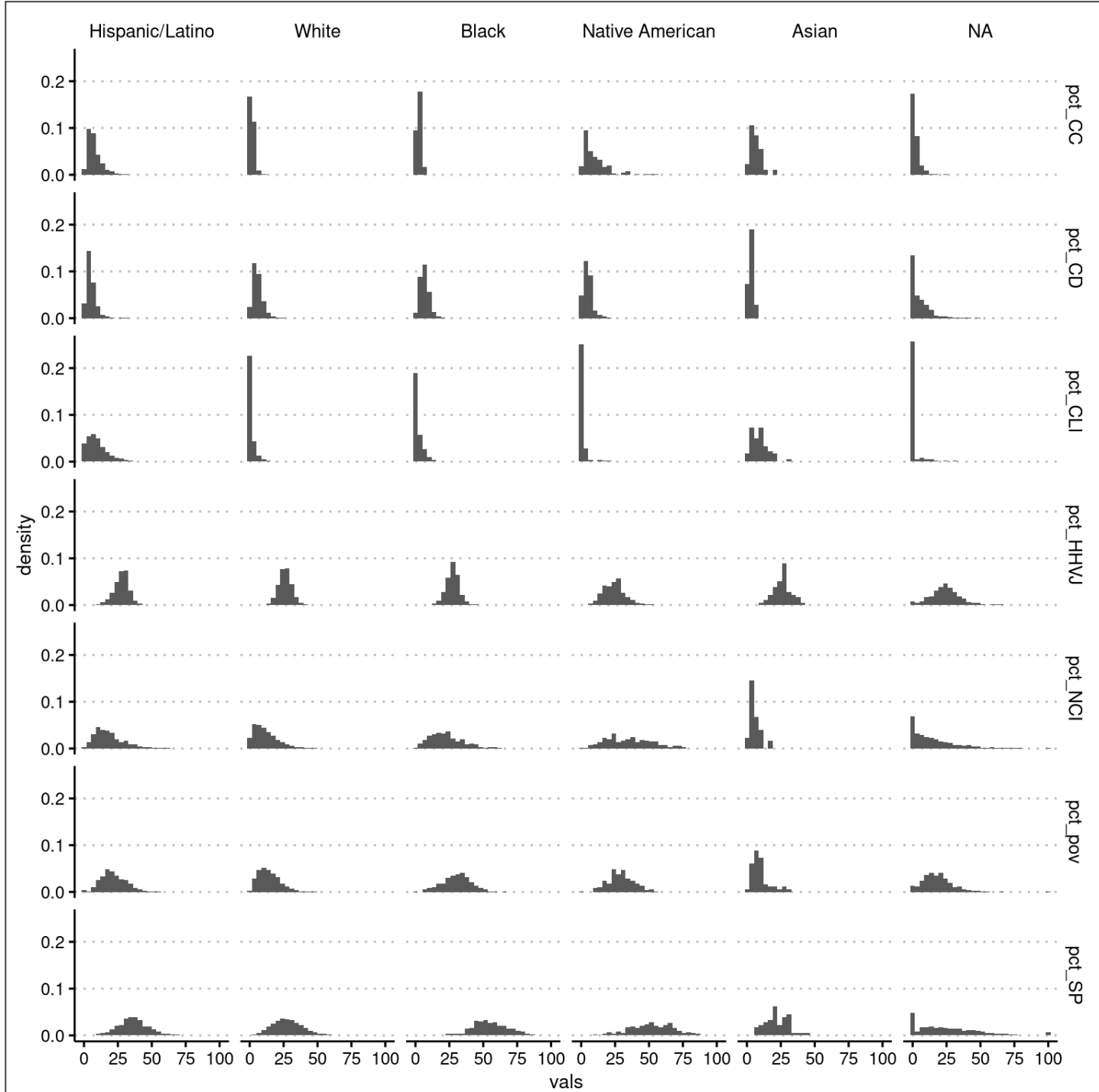# Correlated? | no.

# Regression!

## Linear Regression:

1. Household Condition Data

2. Race Data
   1. %'s of Hispanic/Latino, White, Black, Native American, Asian, and Pacific Islander
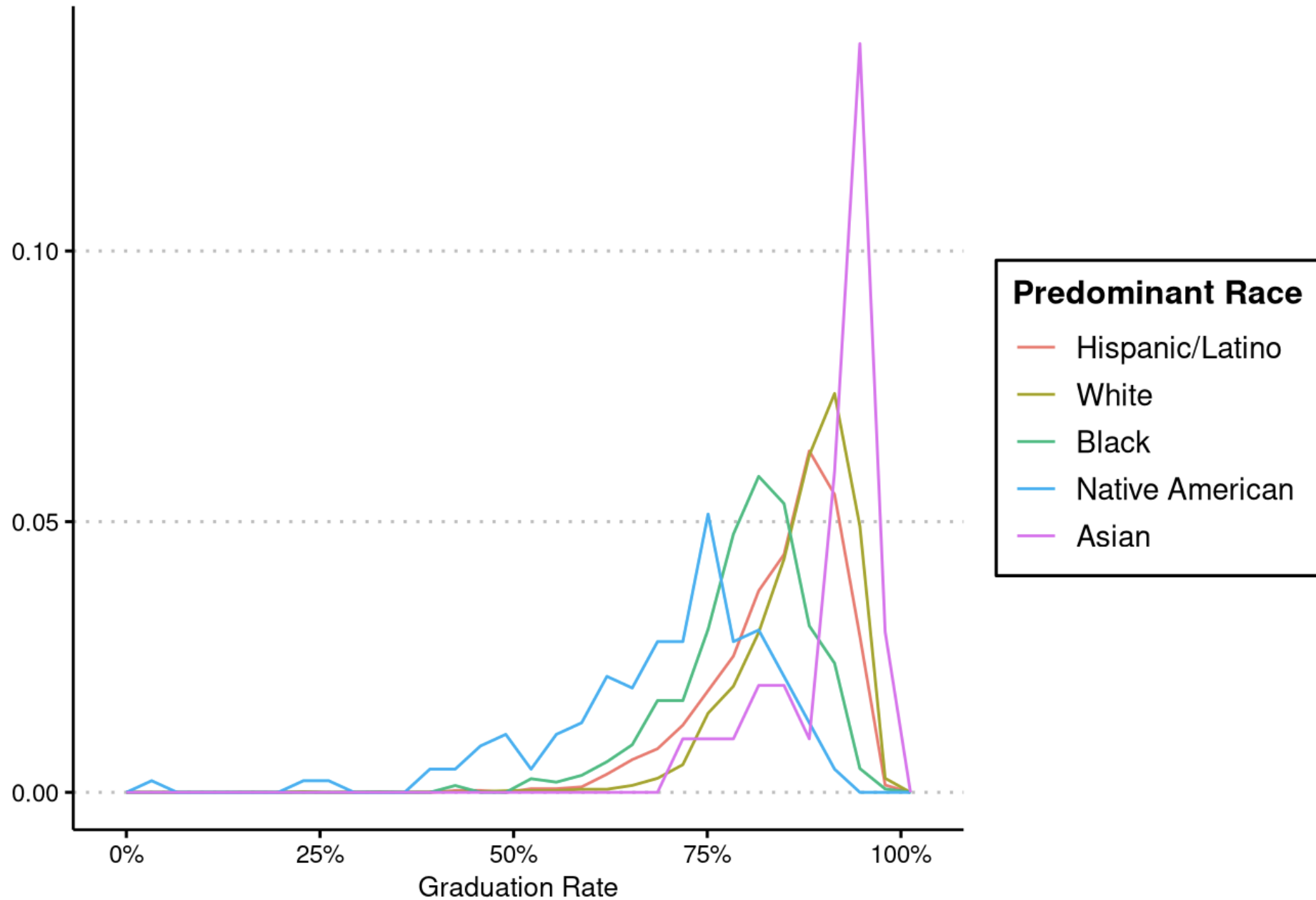
3. Combined

ANOVA to test whether adding race data significantly improves the RSS

Indicator distributions

**Graduation Rate Distribution per Race**

## Regression Moving Forward

- New tool: **Tidy modeling! (tidymodels)**

  - We are trying different regression models beyond just linear models to determine which one makes the most sense.

  - Jon started running models last night which took a while to run.

## Regression Moving Forward

- Preprocessing Steps:

1. Made interaction terms between HH conditions and Racial data

2. Made dummy variables for predominant race

3. Centered and Scaled predictors

4. Removed predictors with near-zero variance

# Models

Tuned parameters with 10-fold cross-validation and random grid search

Linear Regression {lm}

Lasso Regression {glmnet}

Multivariate Adaptive Regression Spline {earth}
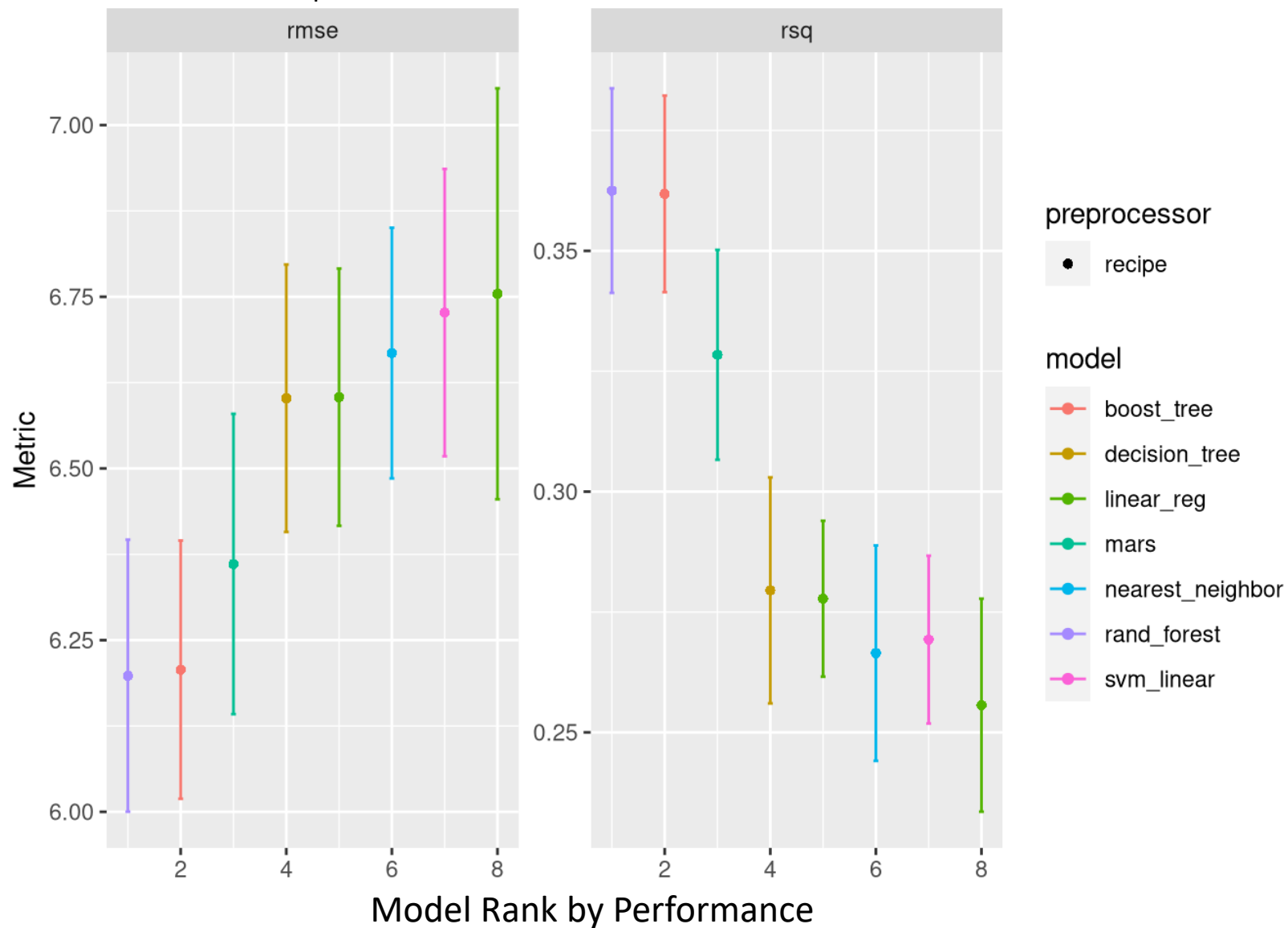
Support Vector Regression {kernlab}

Decision Tree {rpart}

Random Forest {ranger}

Gradient Boosted Trees {xgboost}

K-Nearest Neighbors {kknn}

Model Results:
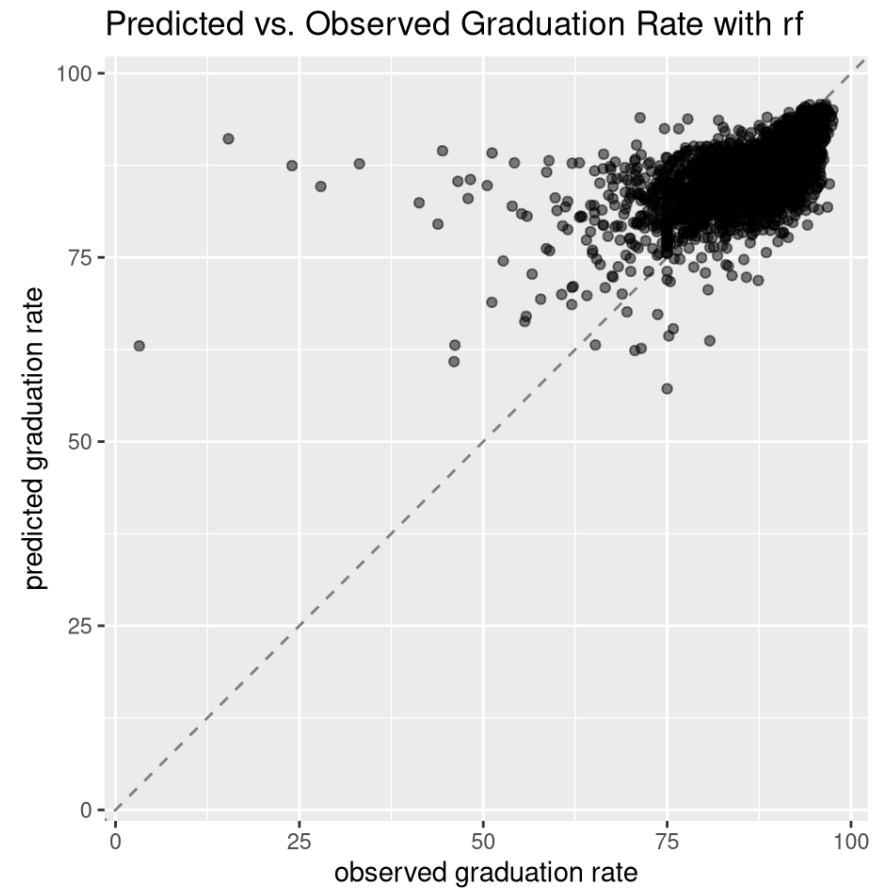
- Test RMSE: 6.77
- Test $R^2$: 0.319

Model Parameters:

- mtry*: 16
- min_n**: 39

*mtry: The number of predictors that will be randomly sampled at each split when creating the tree models.

**min_n: The minimum number of data points in a node that are required for the node to be split further.



Predicted vs. Observed Graduation Rate with rf

# Next Steps…

Add variables to regression model (include financial information)

Determine which type of model is most appropriate – change by region, race

Measure most impactful variables for regression (Societal Implications)

Assessment data