

An aerial photograph of a suburban neighborhood, showing a dense collection of houses with brown roofs and brick walls. The houses are arranged in a grid-like pattern, with some larger houses and some smaller ones. The overall tone is slightly muted, with a dark overlay.

# Household Conditions by Geographic School District

Jon Geiger, Noel Goodwin, Abigail Joppa

Week 9/10



# 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 in the U.S.?

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



## Last Week's Next Steps

Include assessment data in the model

Decide how to split our analysis by region

Interpret variable importance

Evaluate modeling techniques


A large orange circle is positioned on the left side of the slide, partially cut off by the edge.

# Model Variables

## Current Variables:

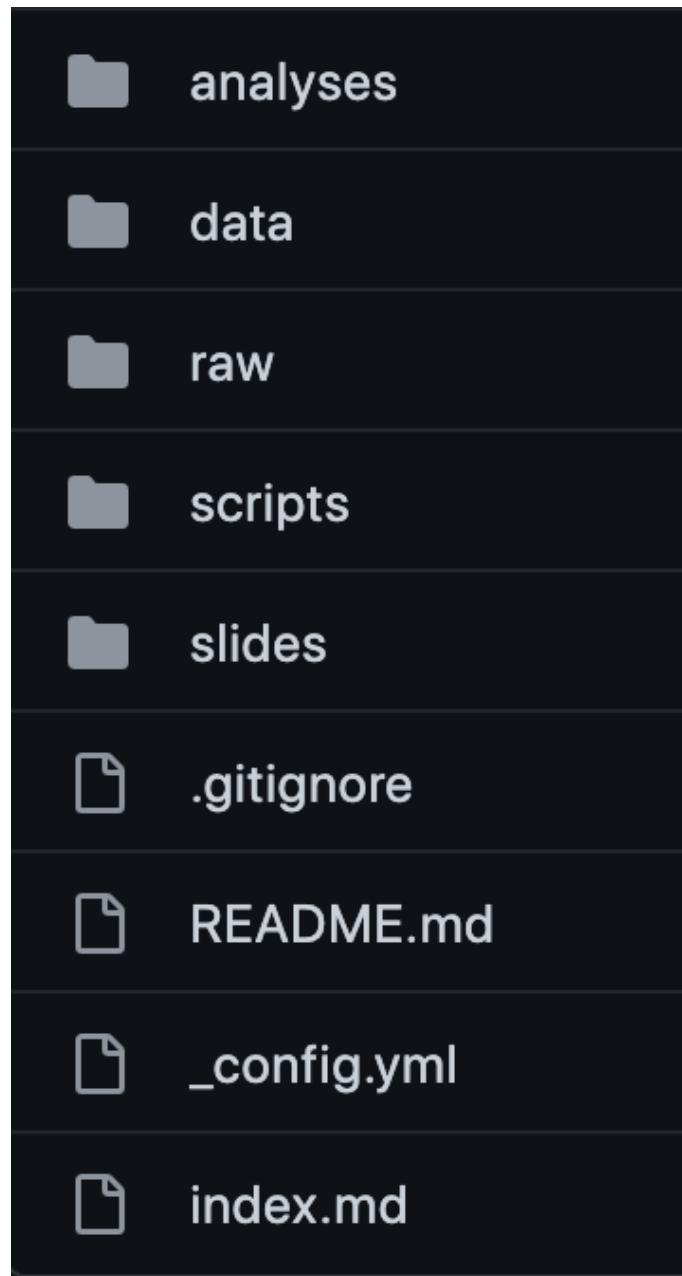
- Household Conditions
- Racial/Ethnic Distribution
- Finance Data
- Assessment Data



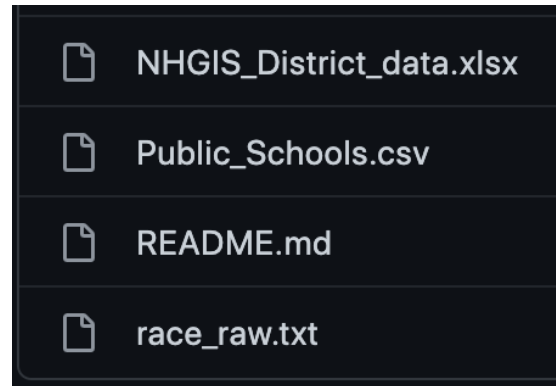


Reproducibility!

# Our Repo



# Our New Data Directories

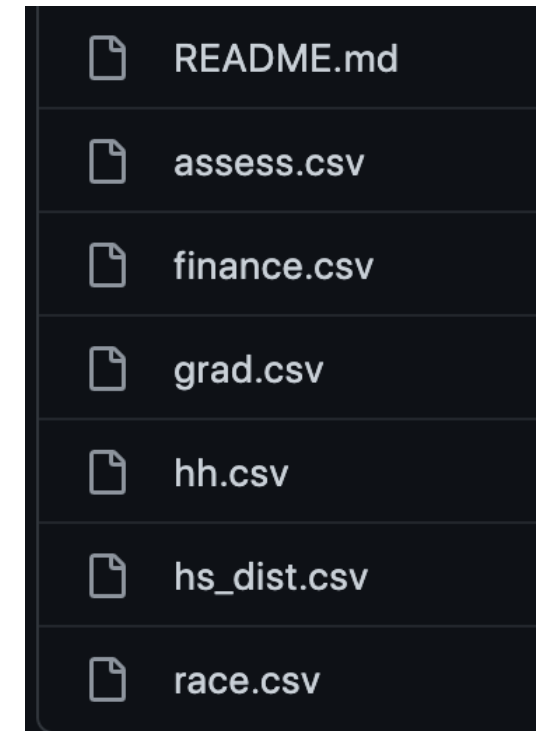


## Our raw data sets

- Original data (NHGIS\_District\_data.xlsx)
- Info about Schools (Public\_Schools.csv)
- Original downloaded race data (race\_raw.txt)


## Our cleaned data sets :


- HH Conditions (hh.csv)
- Grad Rates (grad.csv)
- Race Distribution (race.csv)
- Assessments (assess.csv)
- Financial (finance.csv)
- High School IDs (hs\_dist.csv)






# Scripts

 create\_hh.R


 create\_hs\_dist.R


 create\_race.R


 download\_assess\_data.R

 download\_finance\_data.R

 download\_grad\_data.R

 names\_list.R

 prune\_race\_variables.R

 to\_moe.R



# Data Dictionary

## Household conditions (hh.csv)

Cleaned and pruned version of the original household conditions dataset. Changes include:

- Rename variables for easier use (see [names\_list.R] ([https://github.com/jonmgeiger/household-conditions/blob/main/scripts/names\\_list.R](https://github.com/jonmgeiger/household-conditions/blob/main/scripts/names_list.R)))
- Transform the original "Margin of Error" (MOE) variable into a single sided MOE for future anlayses
- Include LEAID (state fips code + district ID) for easy joining
- Filter out New York Department of Education "School District" -- This is an extreme outlier, as it includes the children for all 32 NY districts. See [create\\_hh.R](#) for more details.

Variables	Description
state	State where each district resides
leaid	Local education agency identification number
dist	District name
children	An estimate of children between the ages 5-17 who are enrolled in school within a certain geographic school district.
pct_pov	Percent of students within each geographical district boundary estimated to be living in poverty
pct_SP	Percent of households within each geographical district boundary estimated to be living in a household with only one father or one mother.

# GitHub Pages

<https://jonmgeiger.github.io/household-conditions/index.html>

## Household Conditions by Geographic School District

Data and Society Capstone Project  
Seattle Pacific University  
By: Jon Geiger, Noel Goodwin, and Abigail Joppa

Home Data Analysis

### Data Dictionary

#### Household Conditions Dataset ([hh.csv](#))

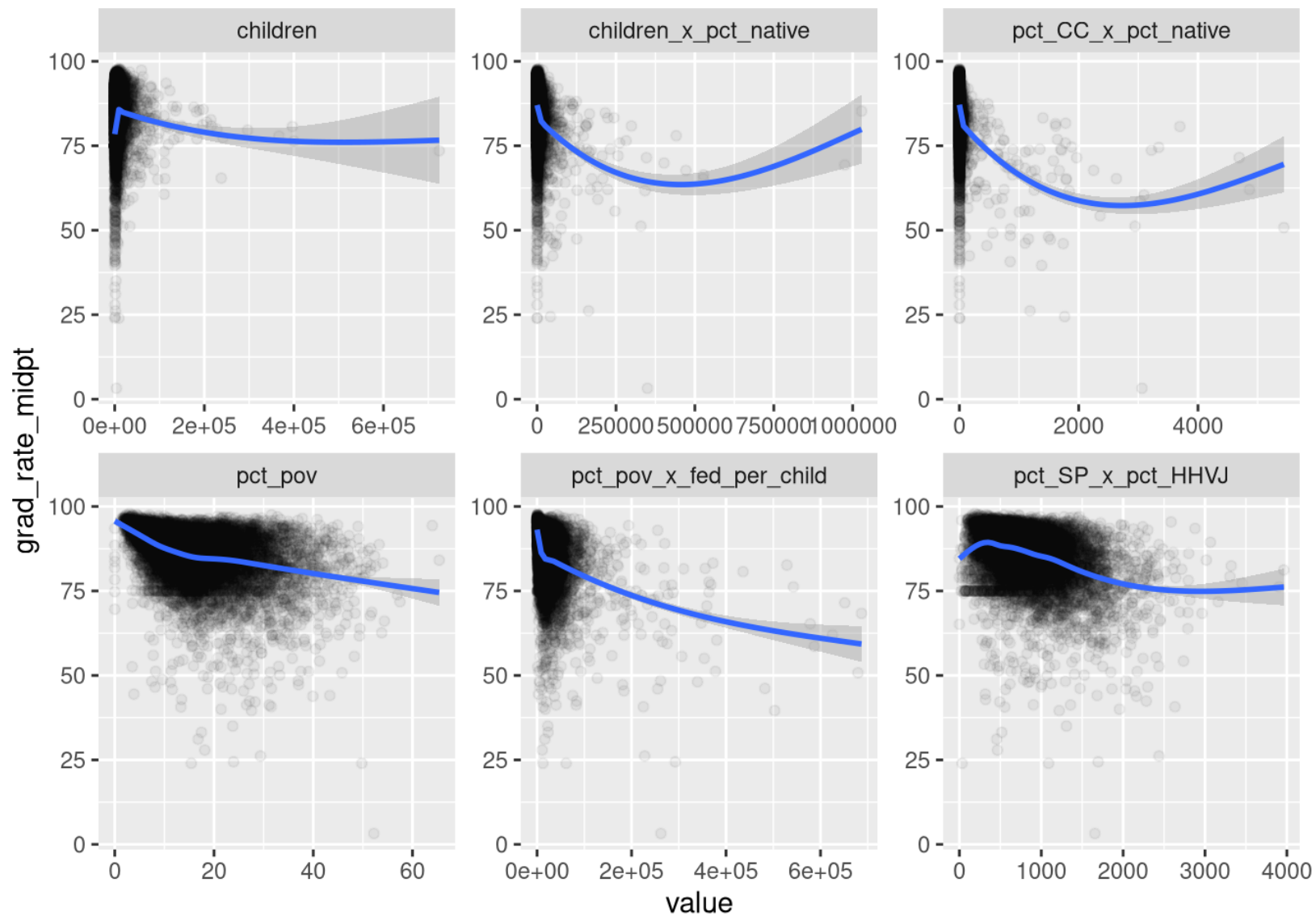
- Provided by the [Urban Institute](#)

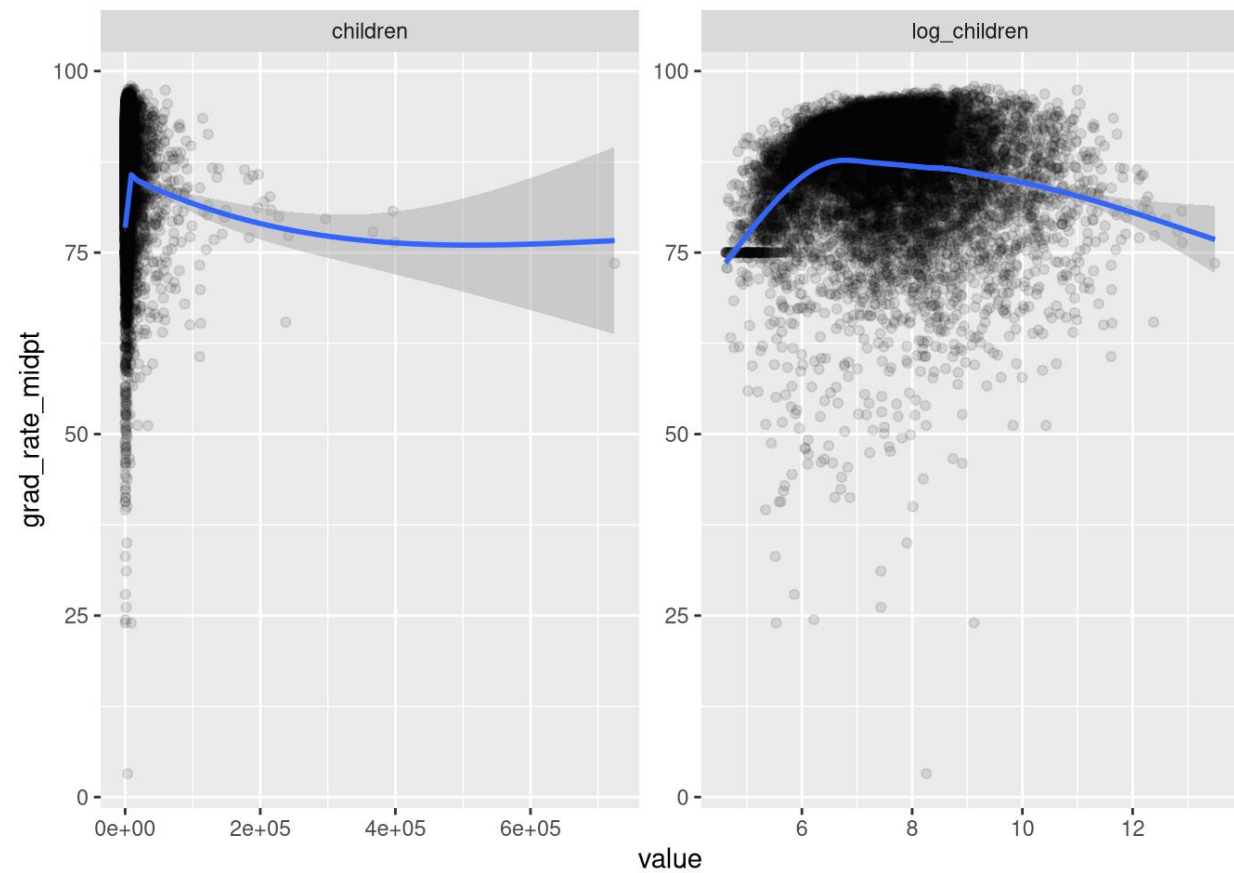
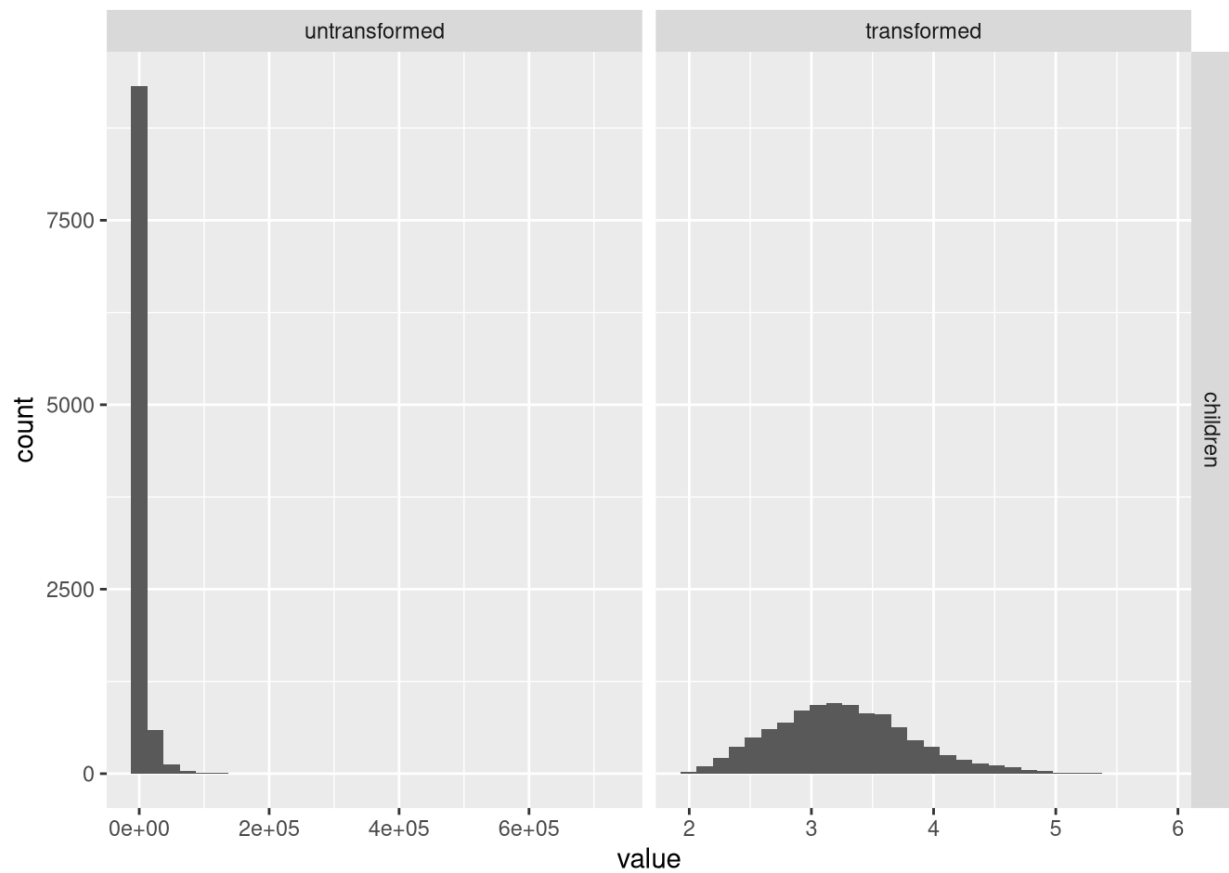
Variables	Description
state	State where each district resides
leaid	Local education agency identification number

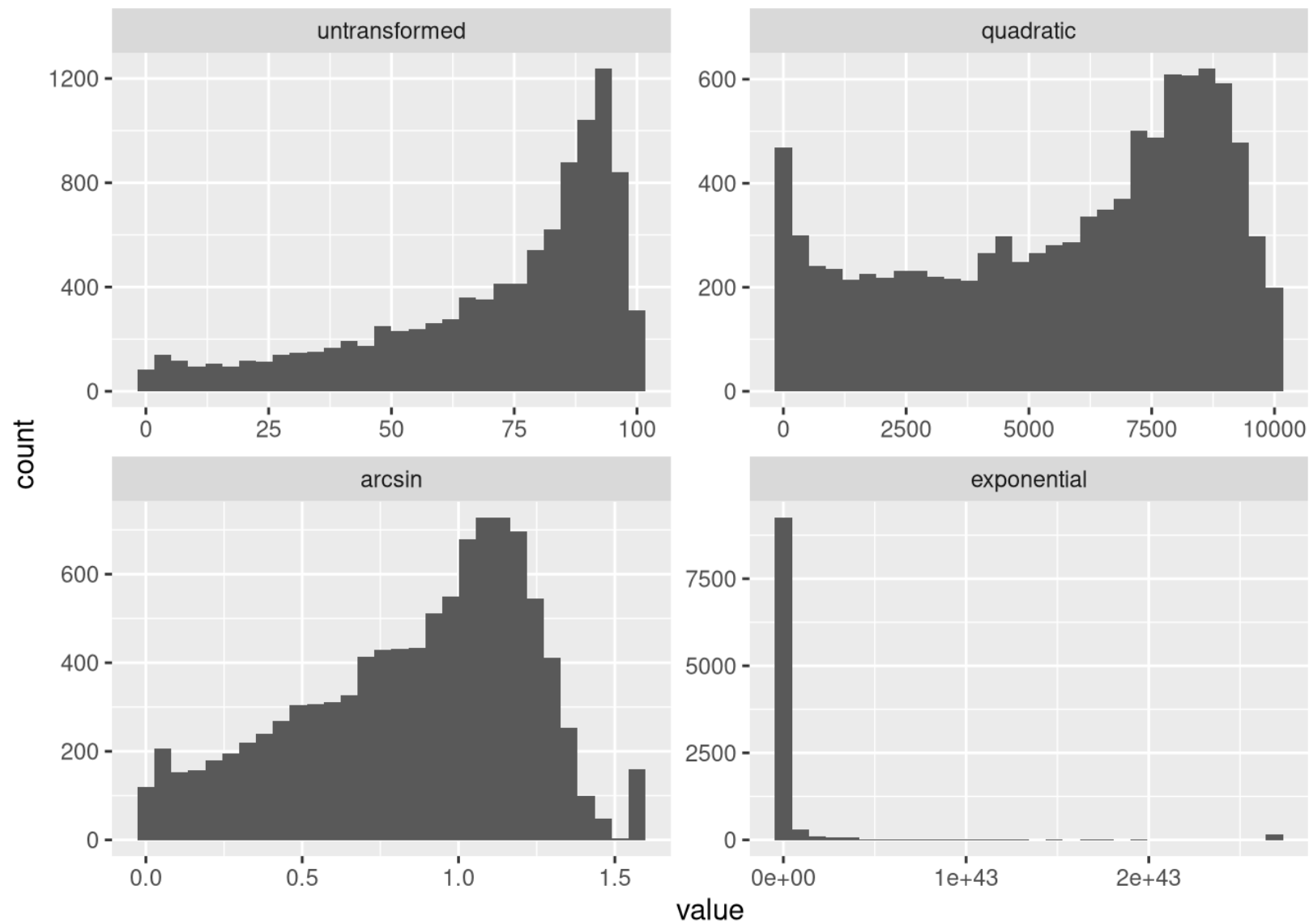


Modeling!

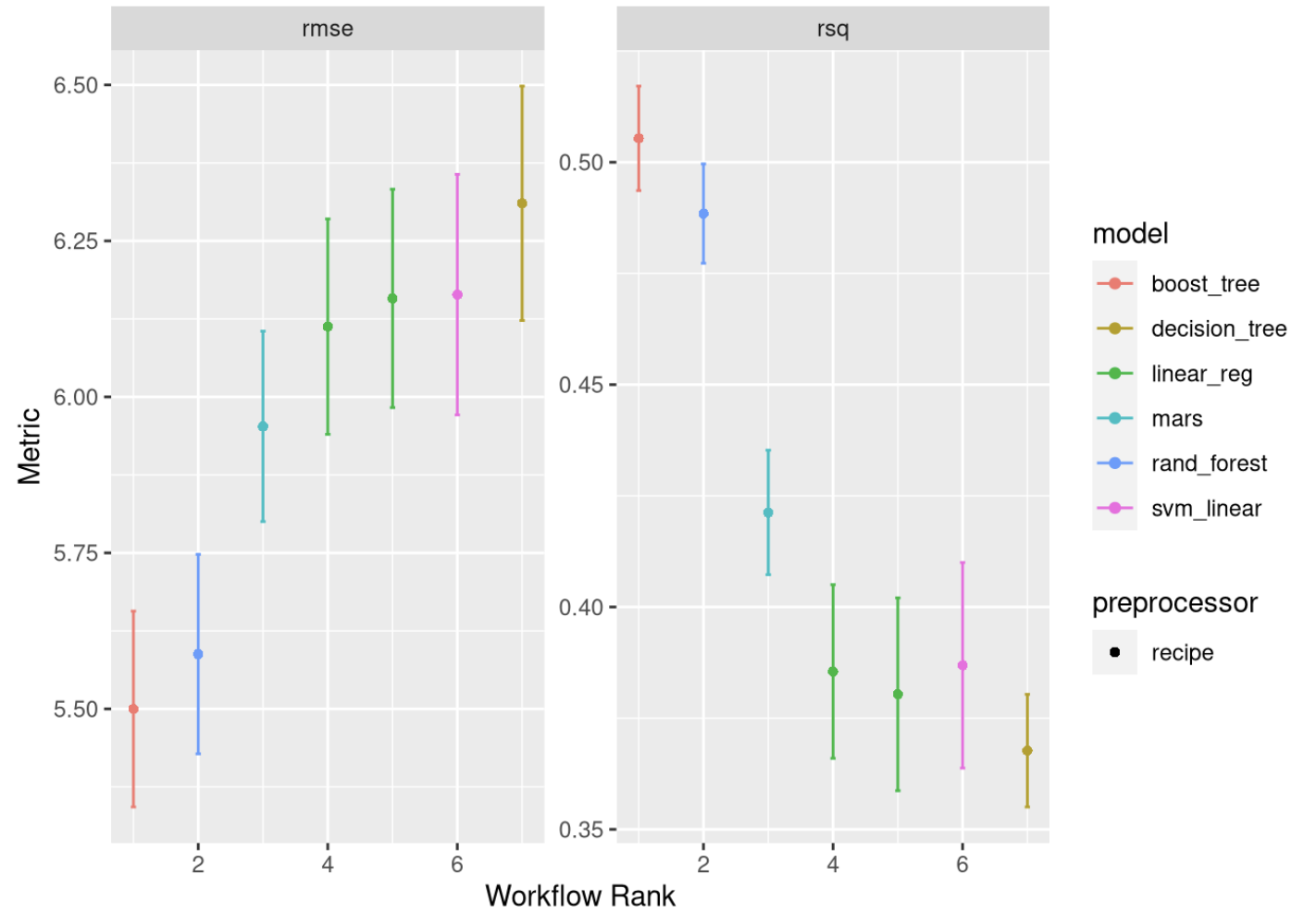
# Investigating Relationships



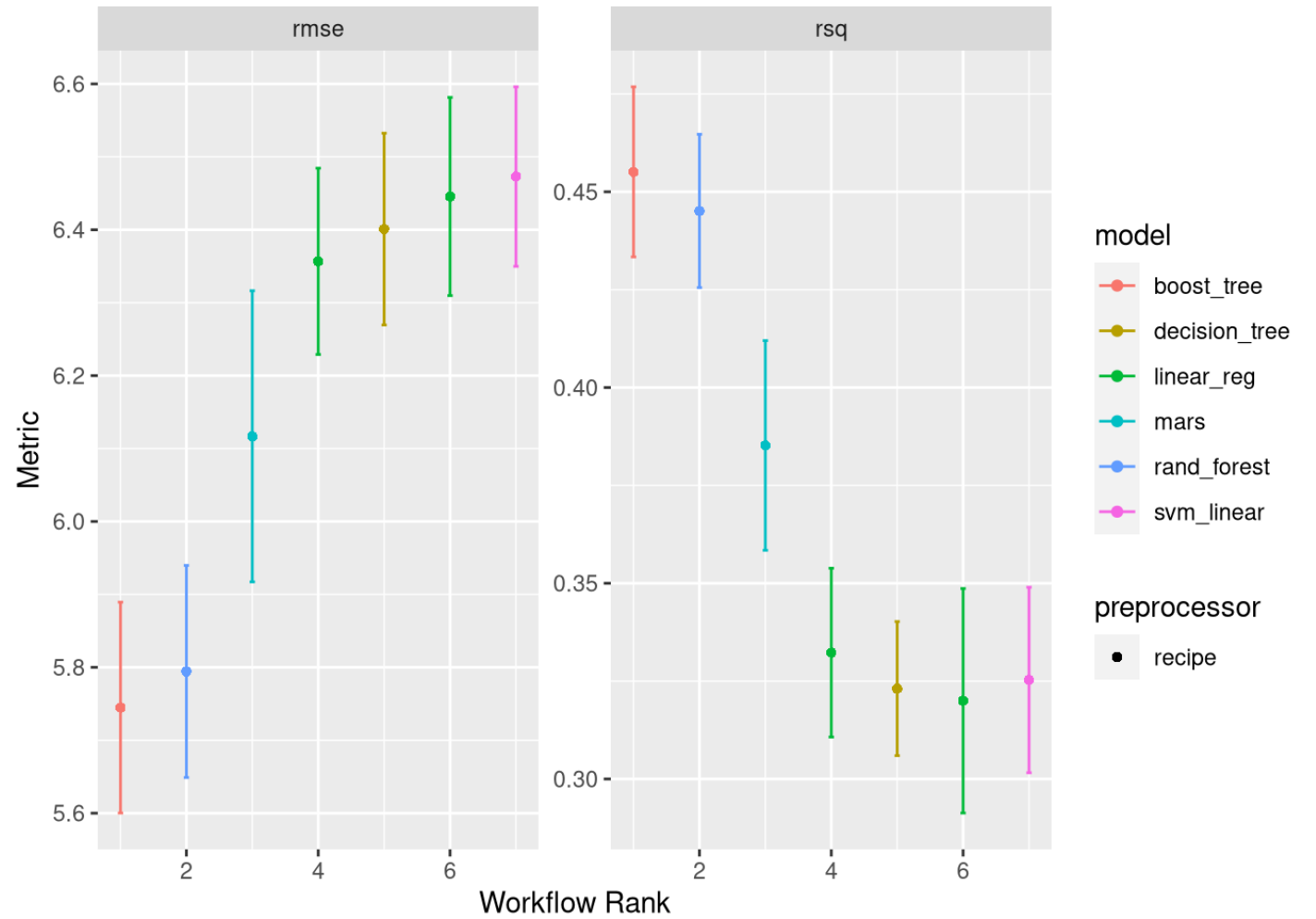




This week's  
results



# Last week's results



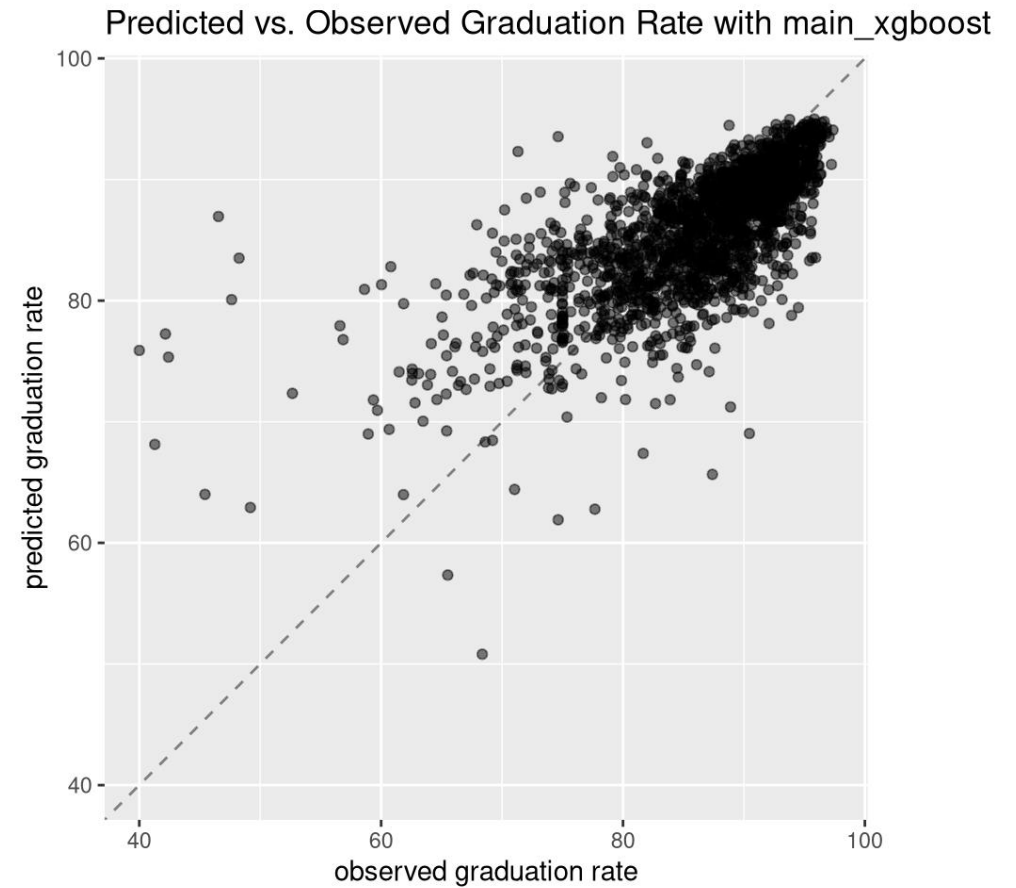


## Result Specifics

```
## # A tibble: 7 × 4
##   wflow_id      rank  rmse   rsq
##   <chr>      <int> <dbl> <dbl>
## 1 main_xgboost      1  5.50 0.505
## 2 main_rf           2  5.59 0.488
## 3 main_mars         3  5.95 0.421
## 4 main_lasso        4  6.11 0.385
## 5 main_lm           5  6.16 0.380
## 6 main_svm          6  6.16 0.387
## 7 main_dtree        7  6.31 0.368
```

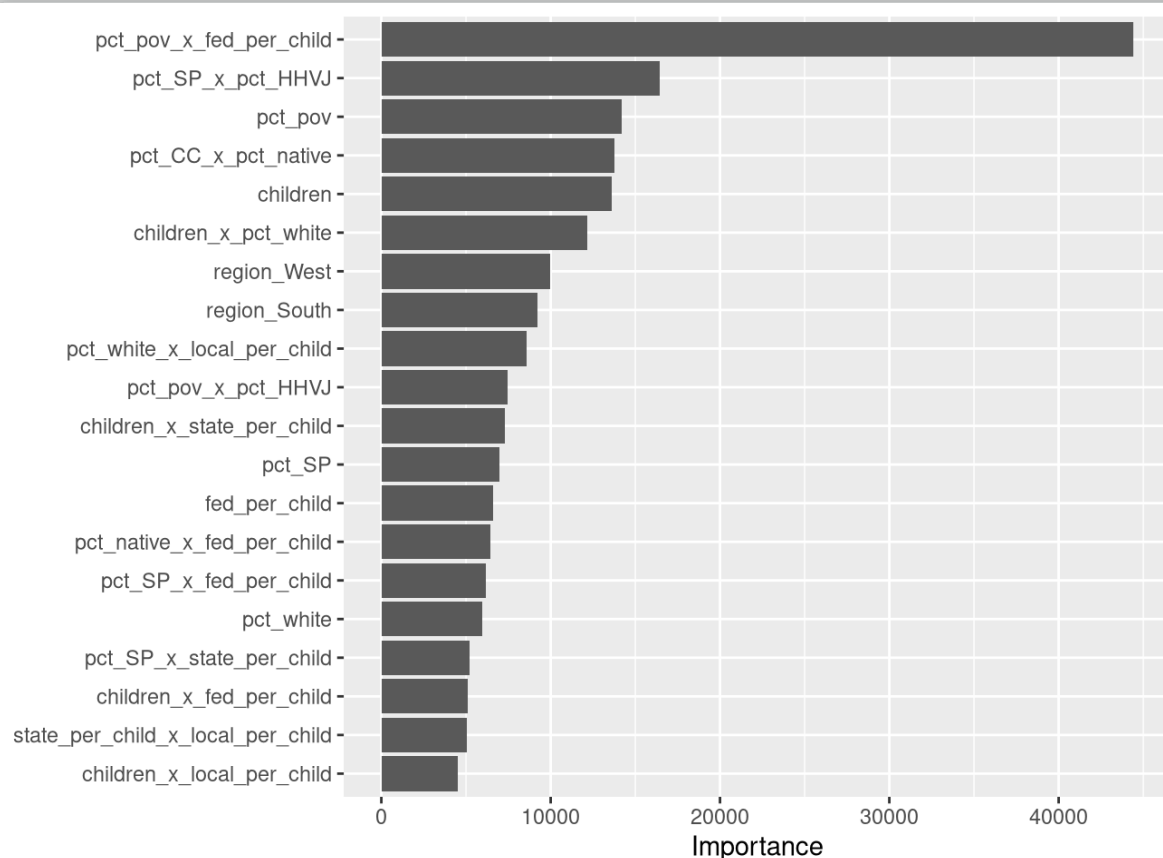
## Model Results:

- Test RMSE: 5.98 --> 5.40
- Test  $R^2$ : 0.451 --> 0.494

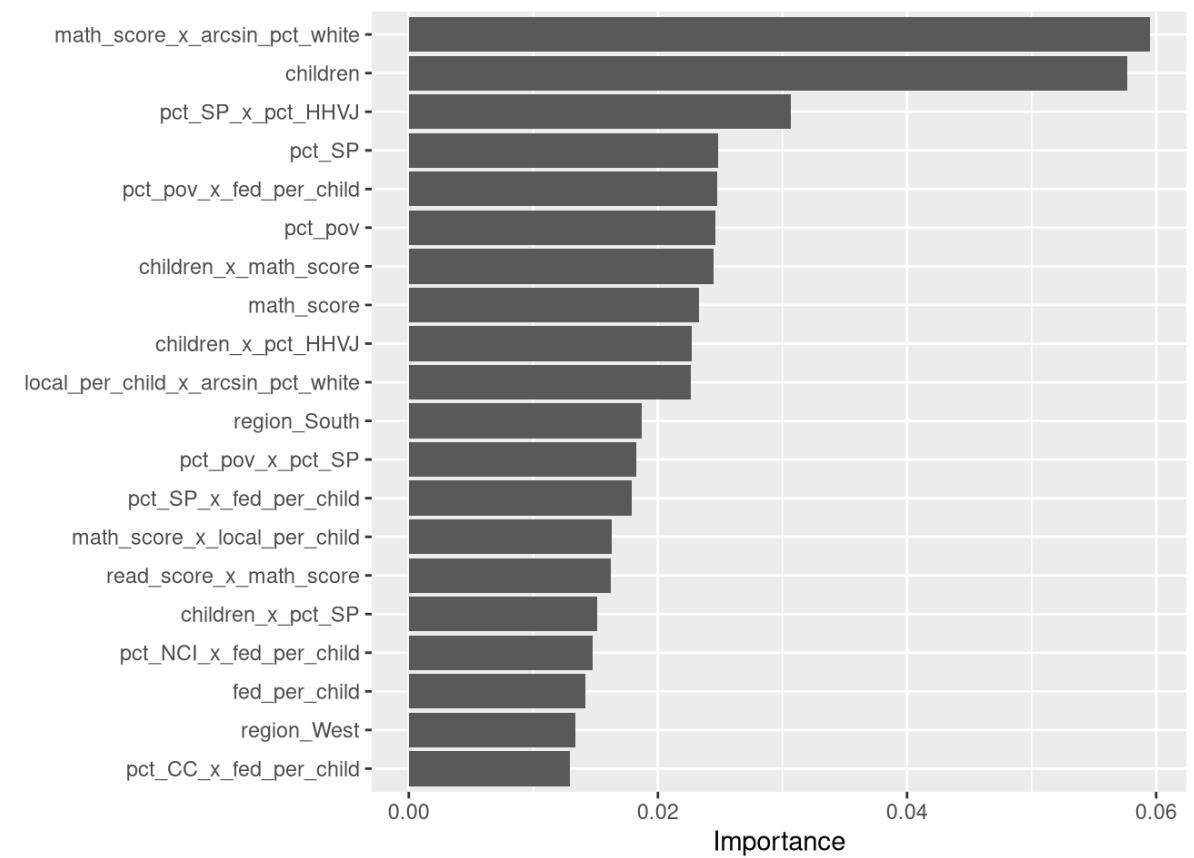


# Variable Importance

Last week's variable importance: Random Forest



This week's variable importance: Boosted Tree



## Next Steps

Work on Final Presentation & Deliverables (Our Data Story)

Include only districts with high schools (interpreting grad rates?)

Compare model accuracy with four regional models OR state-wide dummy variables instead of region