



Overview

Clustering

*Machine
Learning*

*EDA and
Visualization*

*Hypothesis
Testing*

Conclusions

Practicum I

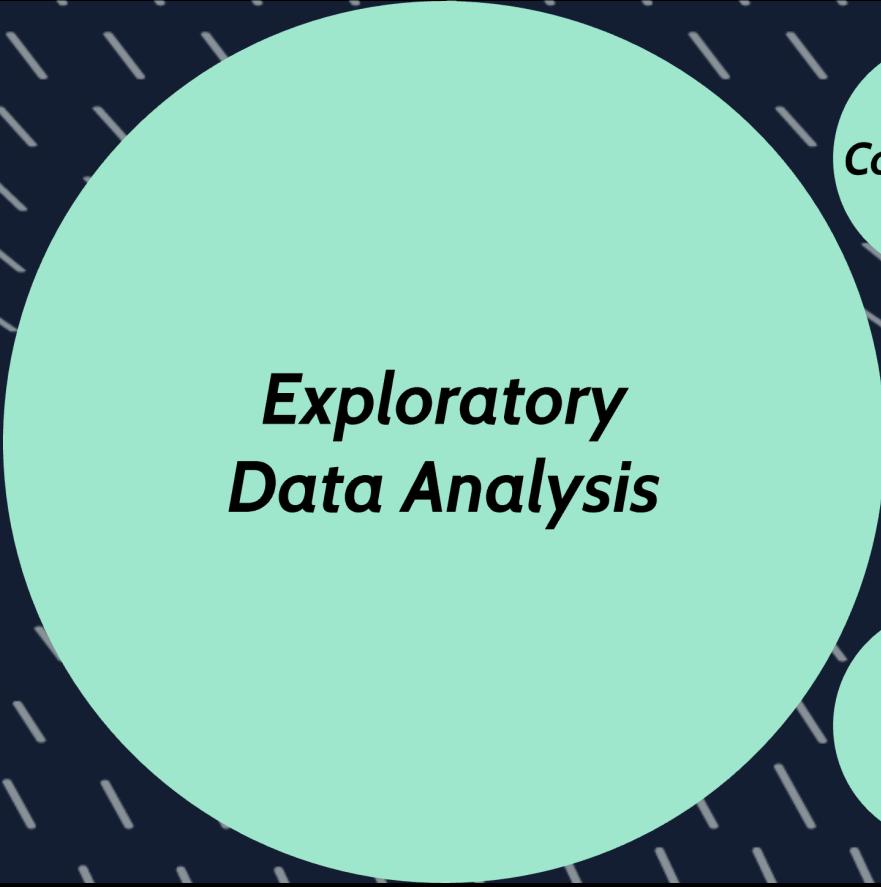
Colorado Department of Education Data Analysis
Lee Wise

Data Preparation

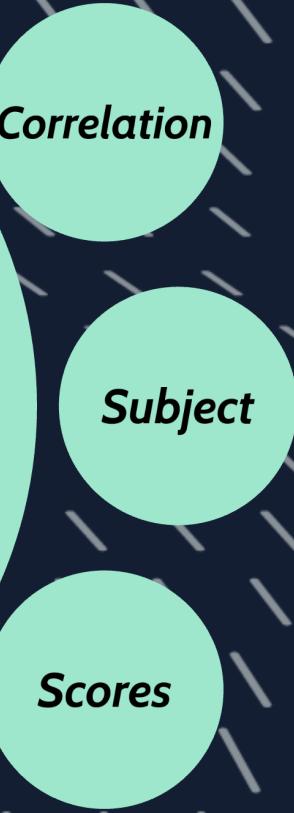
- Data gathered from CDE website for 2016, 2017, 2018.
- Data pulled from IRS website.
- Missing values imputed with average.
- Growth and PWR Data Removed
- Final datasets consist of a school-wide dataset with 56,789 observations and 33 features

EDA

Dashboard



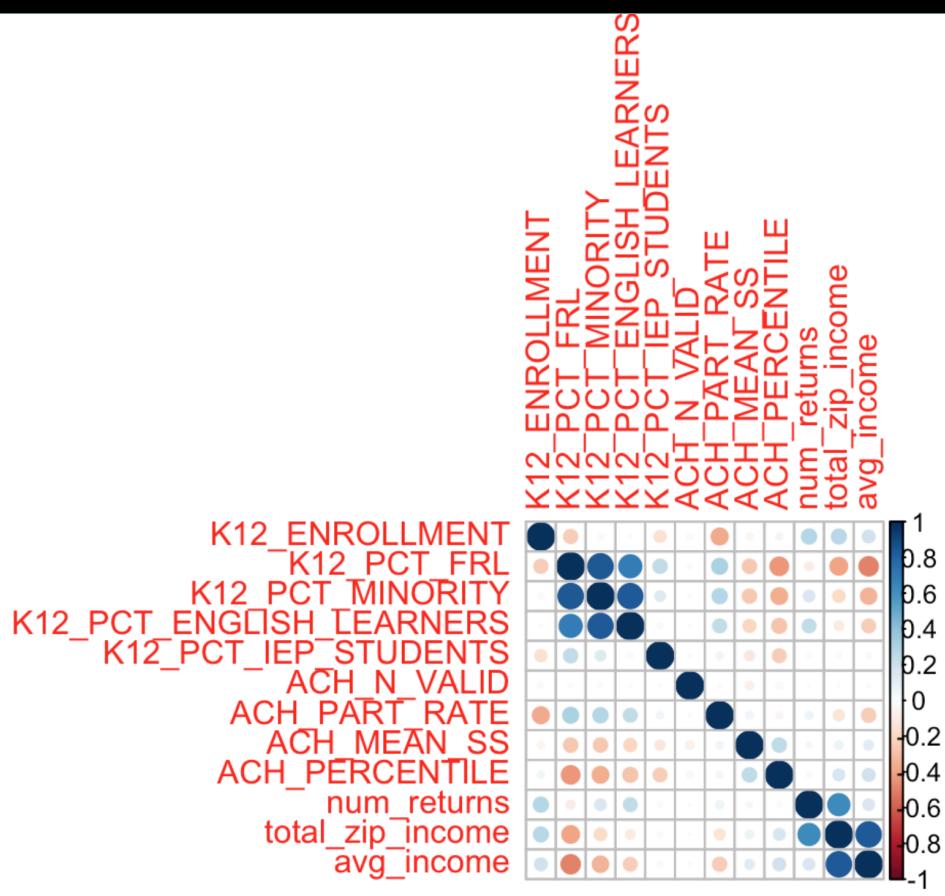
Exploratory Data Analysis

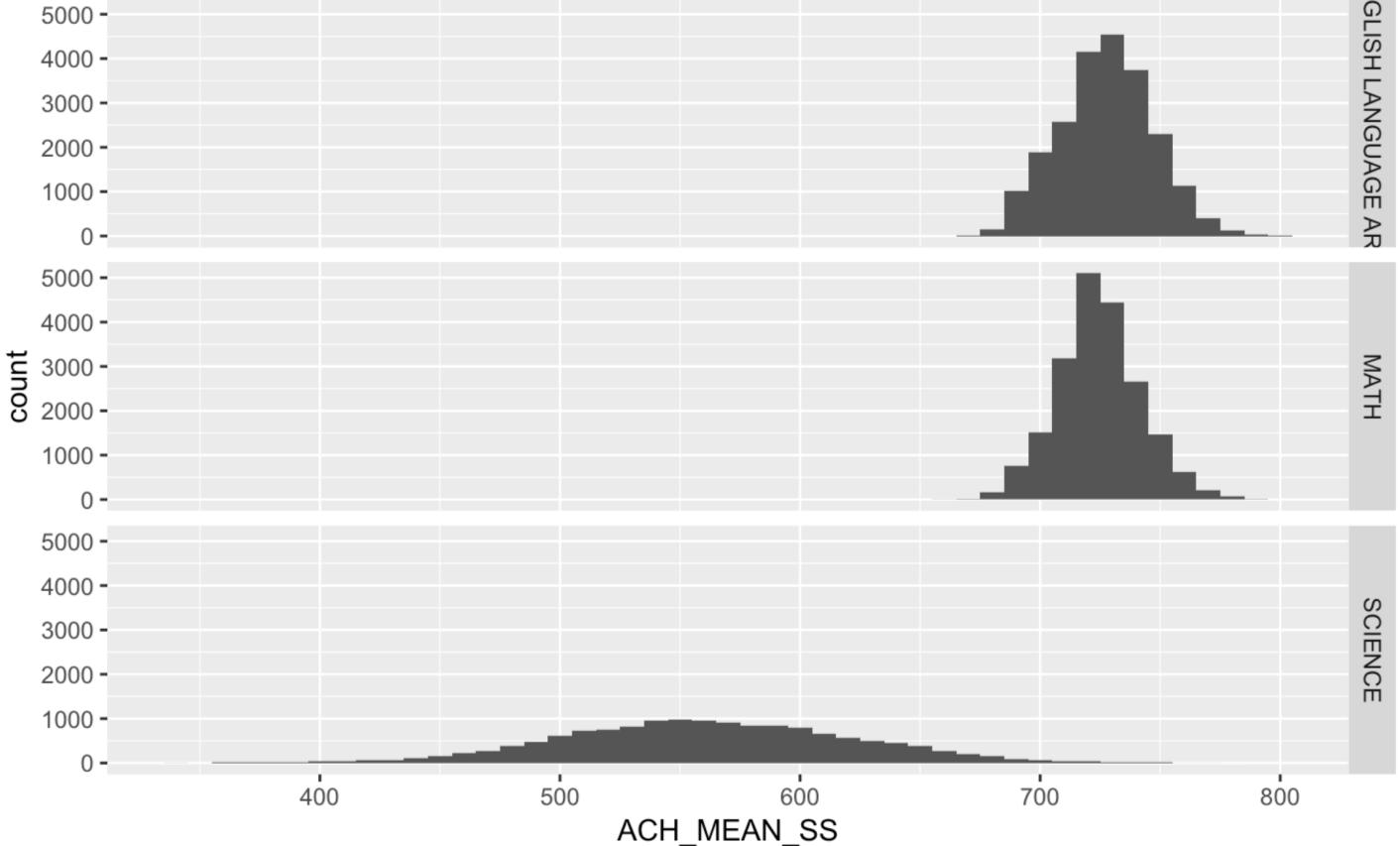


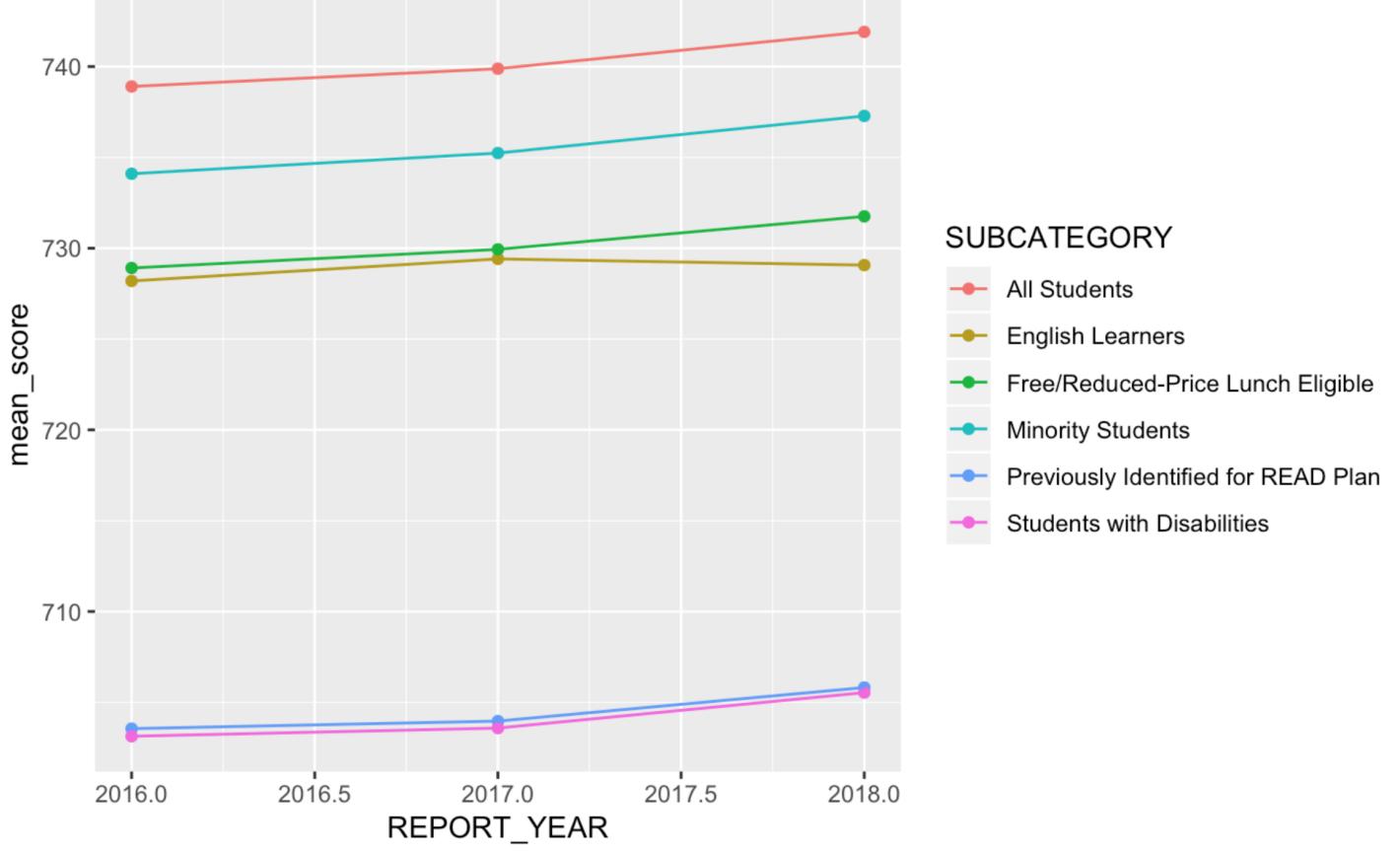
Correlation

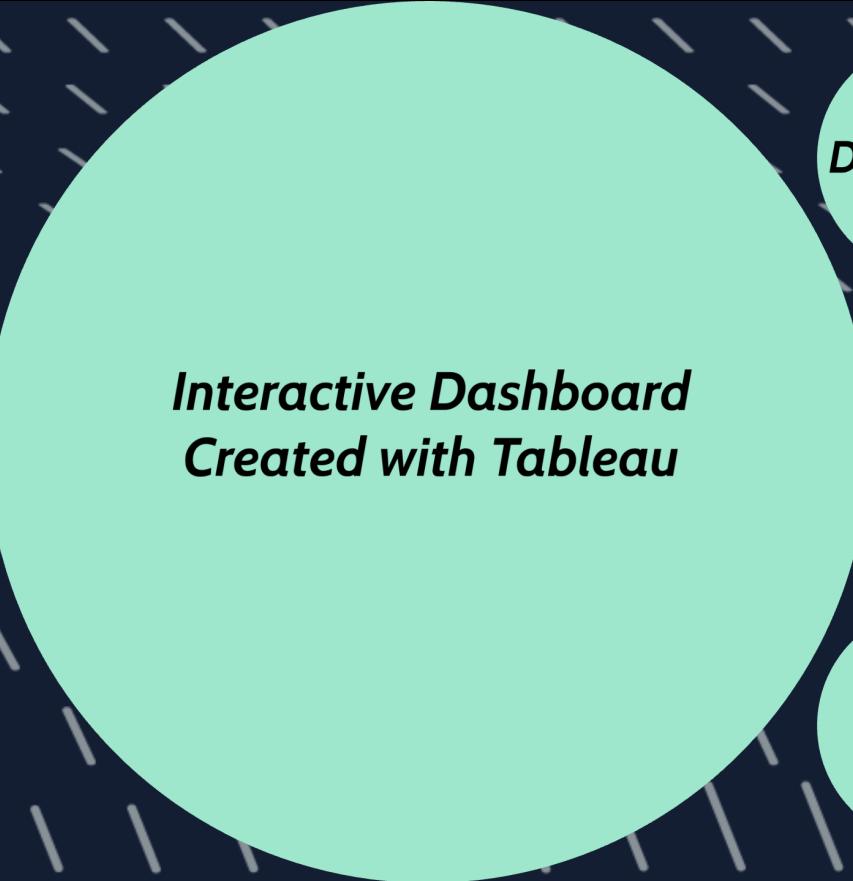
Subject

Scores









***Interactive Dashboard
Created with Tableau***



Dashboard

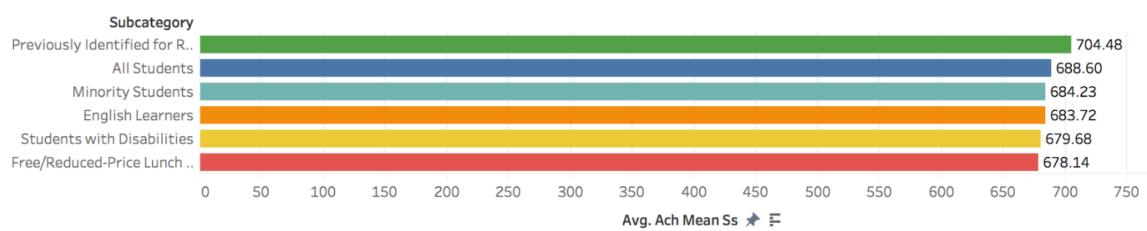


Selecting

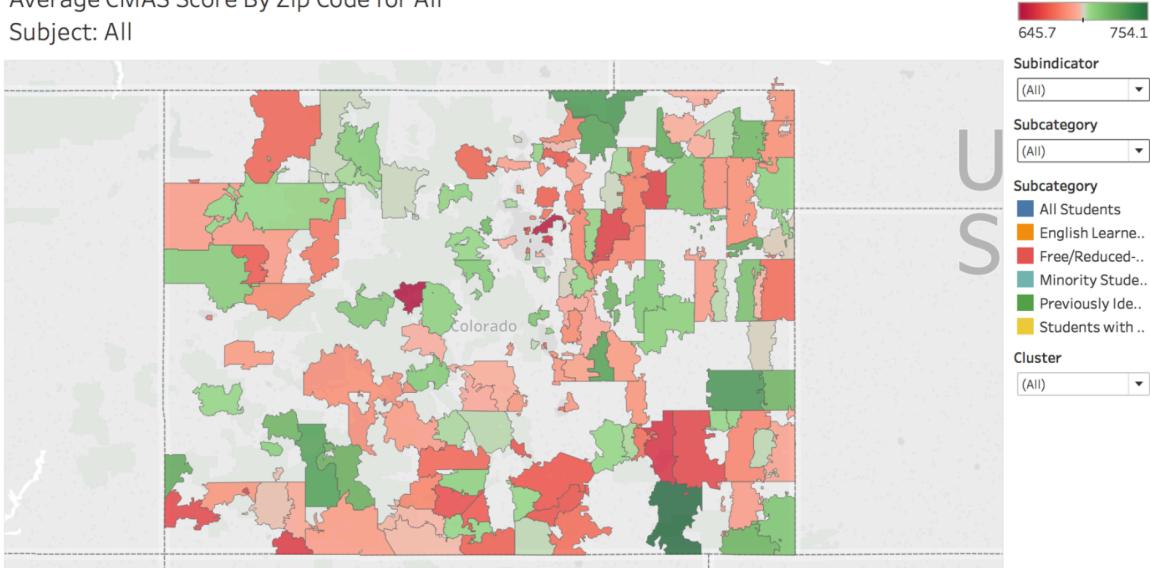


Filtering

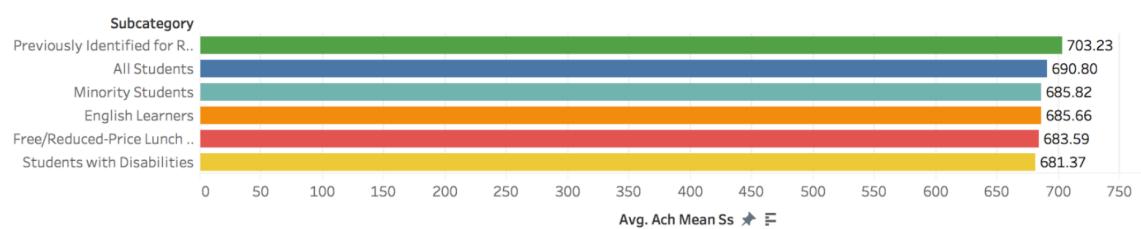
Average CMAS Score



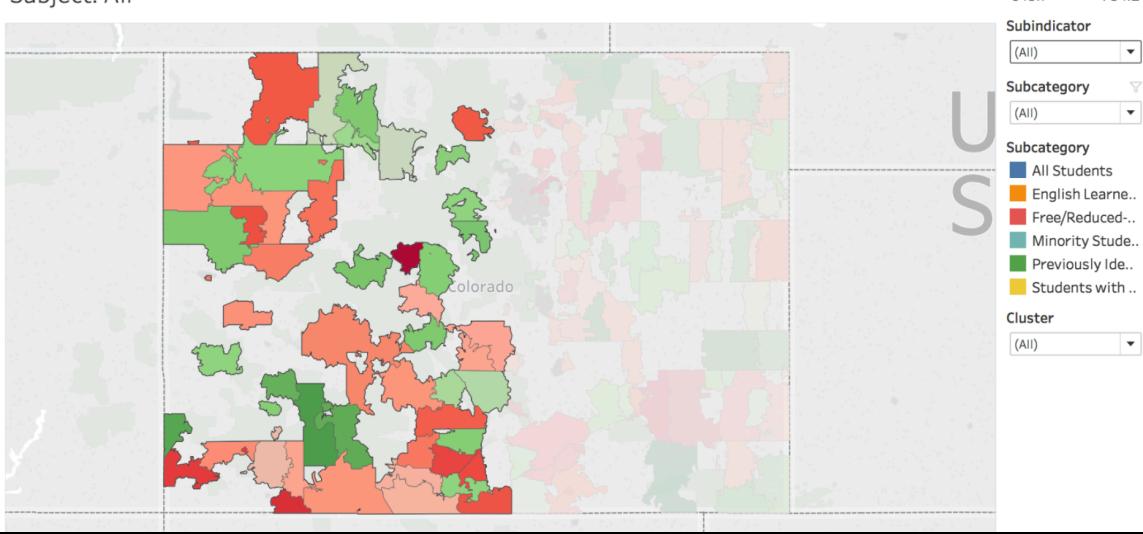
Average CMAS Score By Zip Code for All Subject: All



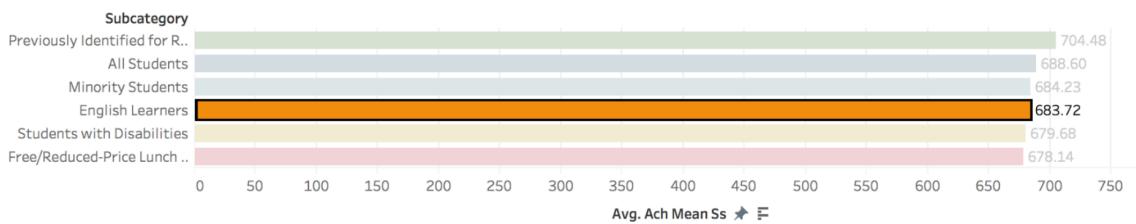
Average CMAS Score



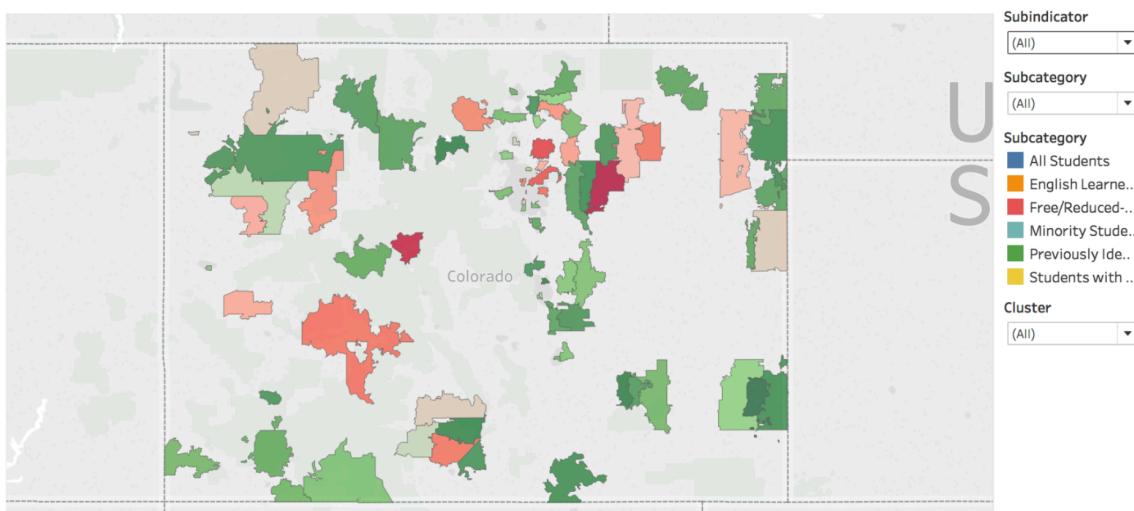
Average CMAS Score By Zip Code for All Subject: All

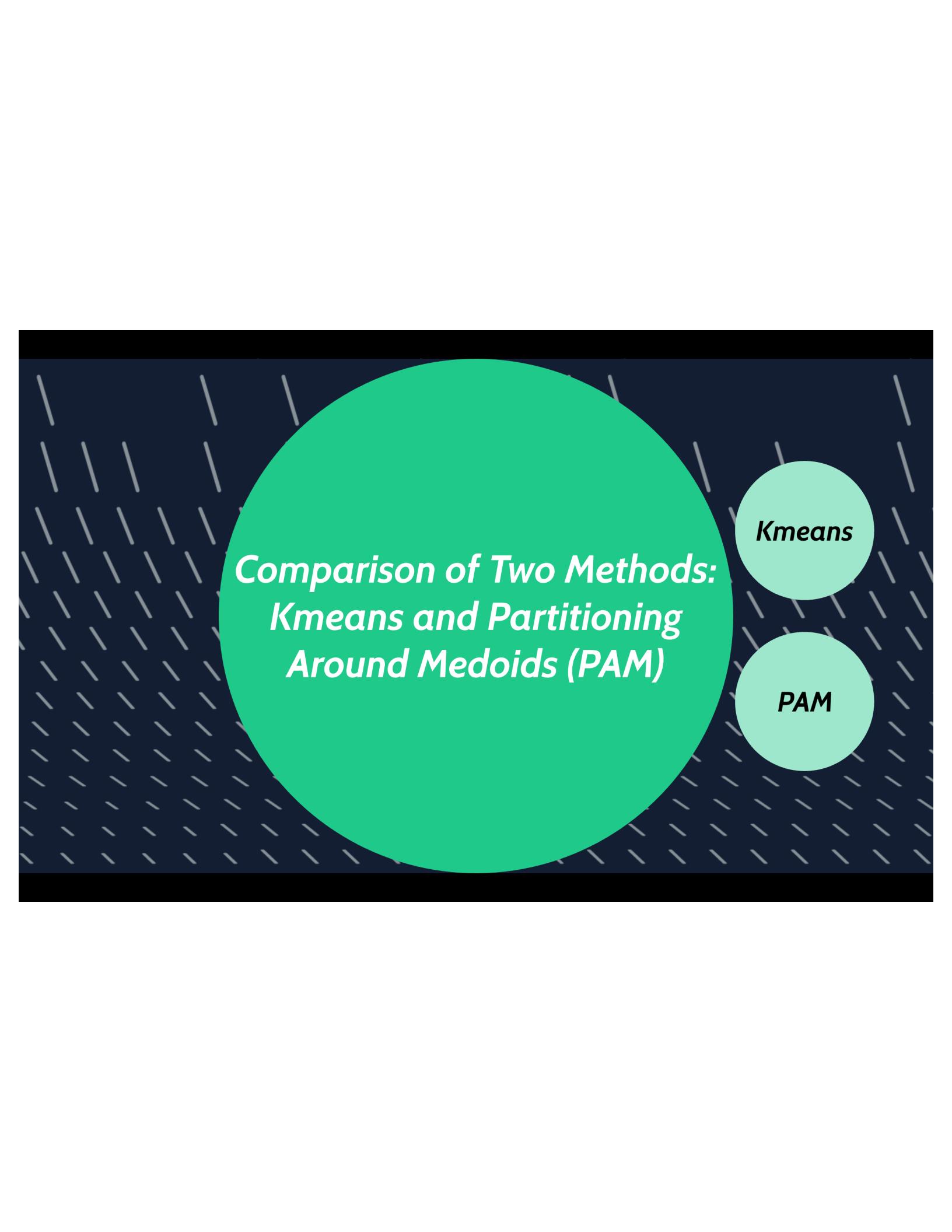


Average CMAS Score



Average CMAS Score By Zip Code for None Subject: All





Comparison of Two Methods: Kmeans and Partitioning Around Medoids (PAM)

Kmeans

PAM

Kmeans

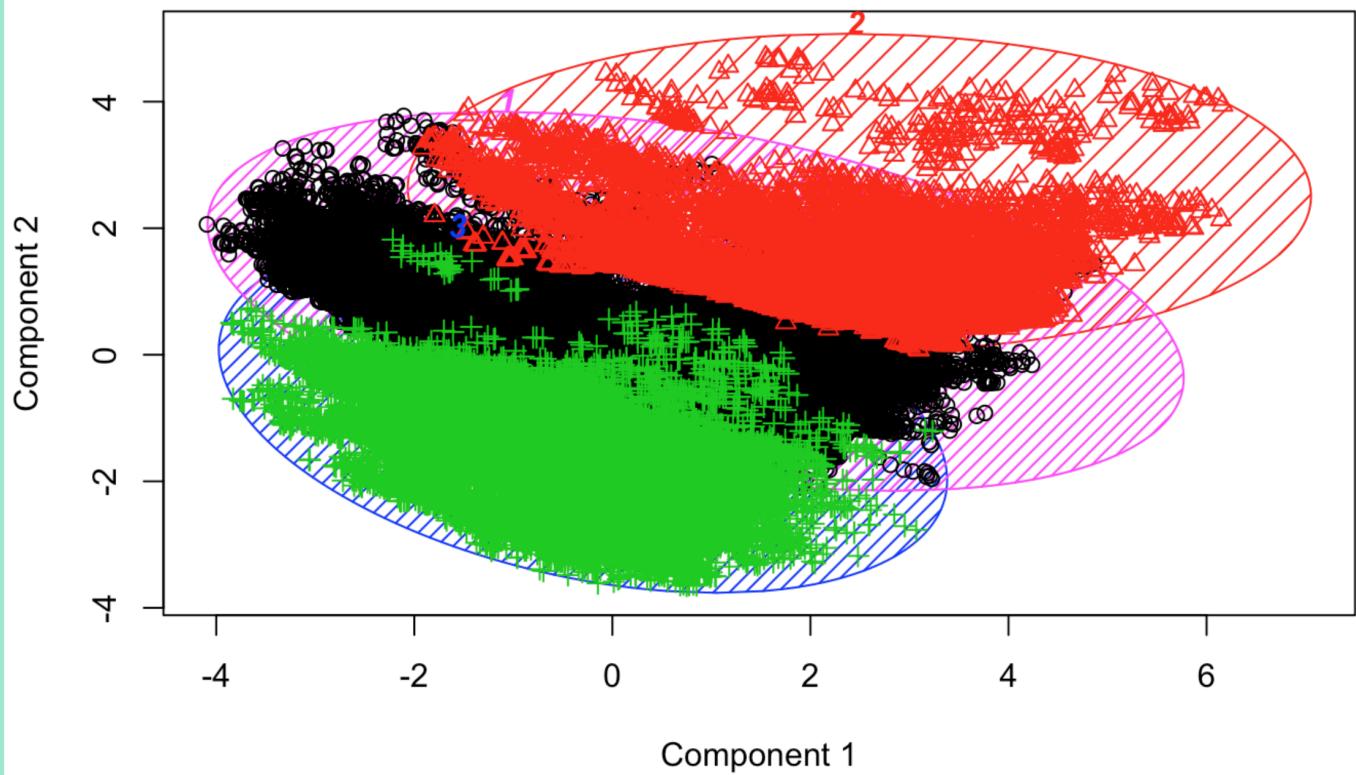
- Moves centroids
- Minimizes squared error

Unscaled

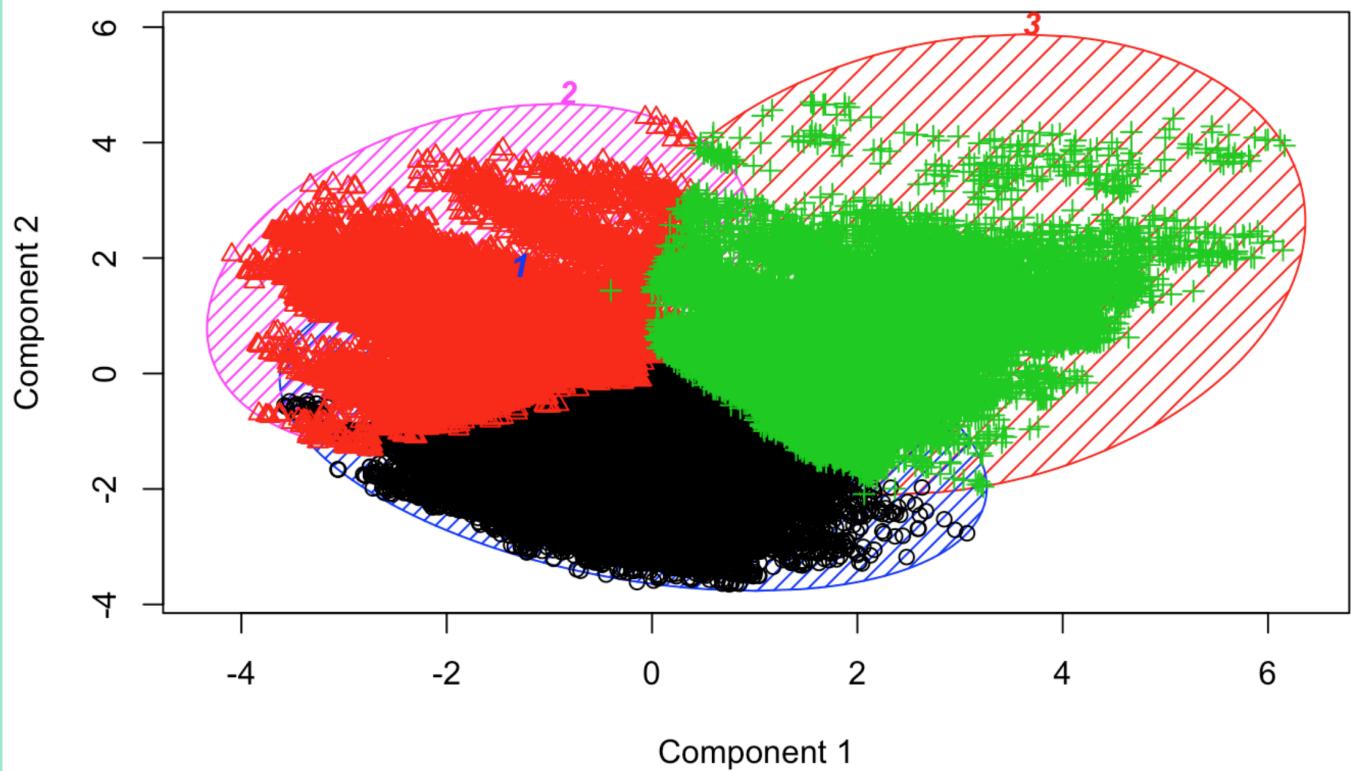
Scaled

Clusters

2D representation of the Cluster solution



2D representation of the Cluster solution



Characteristics

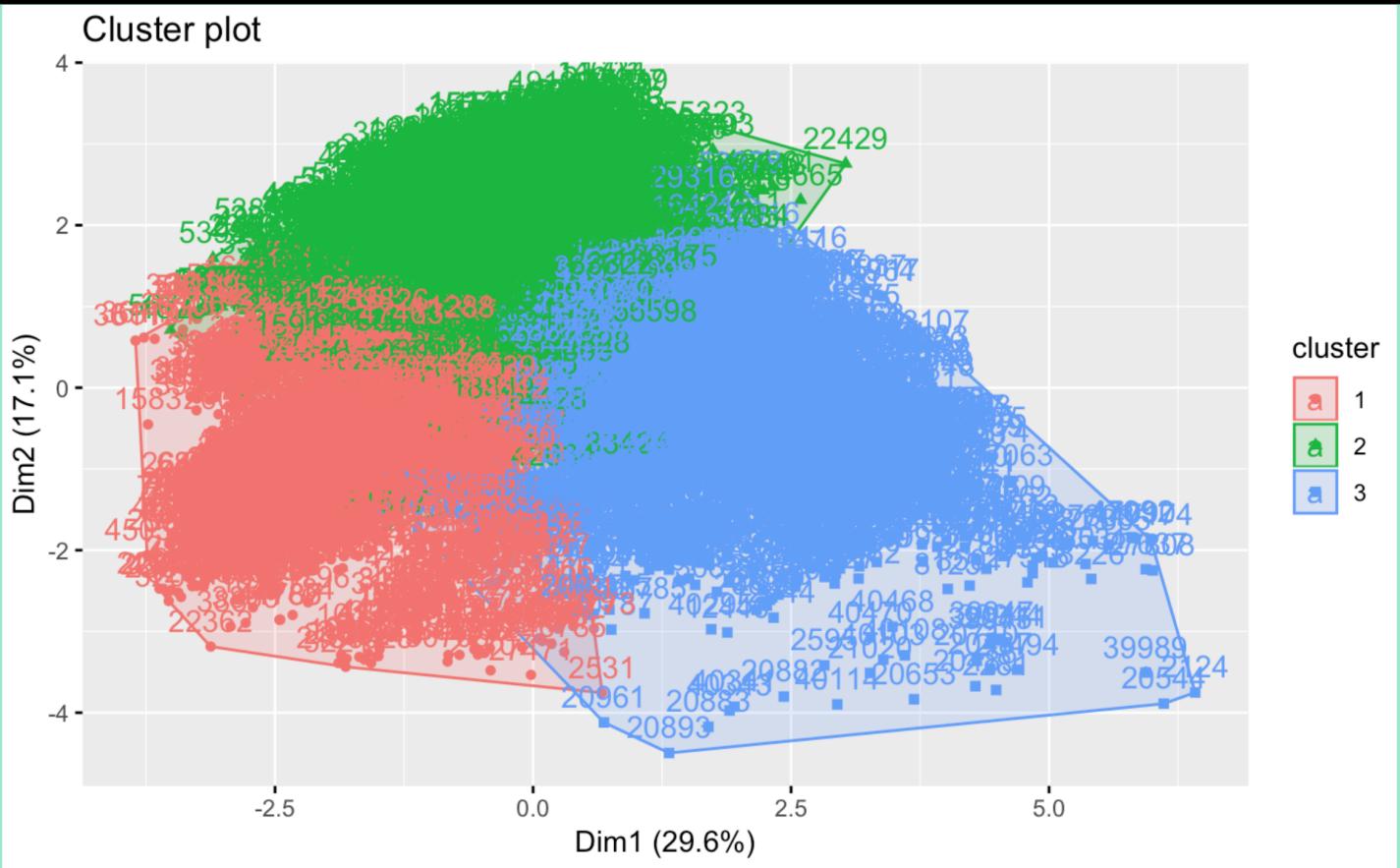
- *Cluster 1: High enrollment, high minority and percent of English learners, lowest test scores.*
- *Cluster 2: High enrollment, low free/reduced lunch/ minority/English learner percentages, highest average income by a significant margin, highest test scores.*
- *Cluster 3: Low enrollment, high percent free/reduced lunch, medium test scores, low income.*

Partitioning Around Medoids (PAM)

- Changes medoids
(data point)
- Minimizes
pairwise distance

Clusters

Comparison



Agreement Between Kmeans and PAM

		Reference		
Prediction	1	2	3	
1	1679	51	33	
2	96	1737	262	
3	6	9	1806	

Overall Statistics

Accuracy : 0.9195

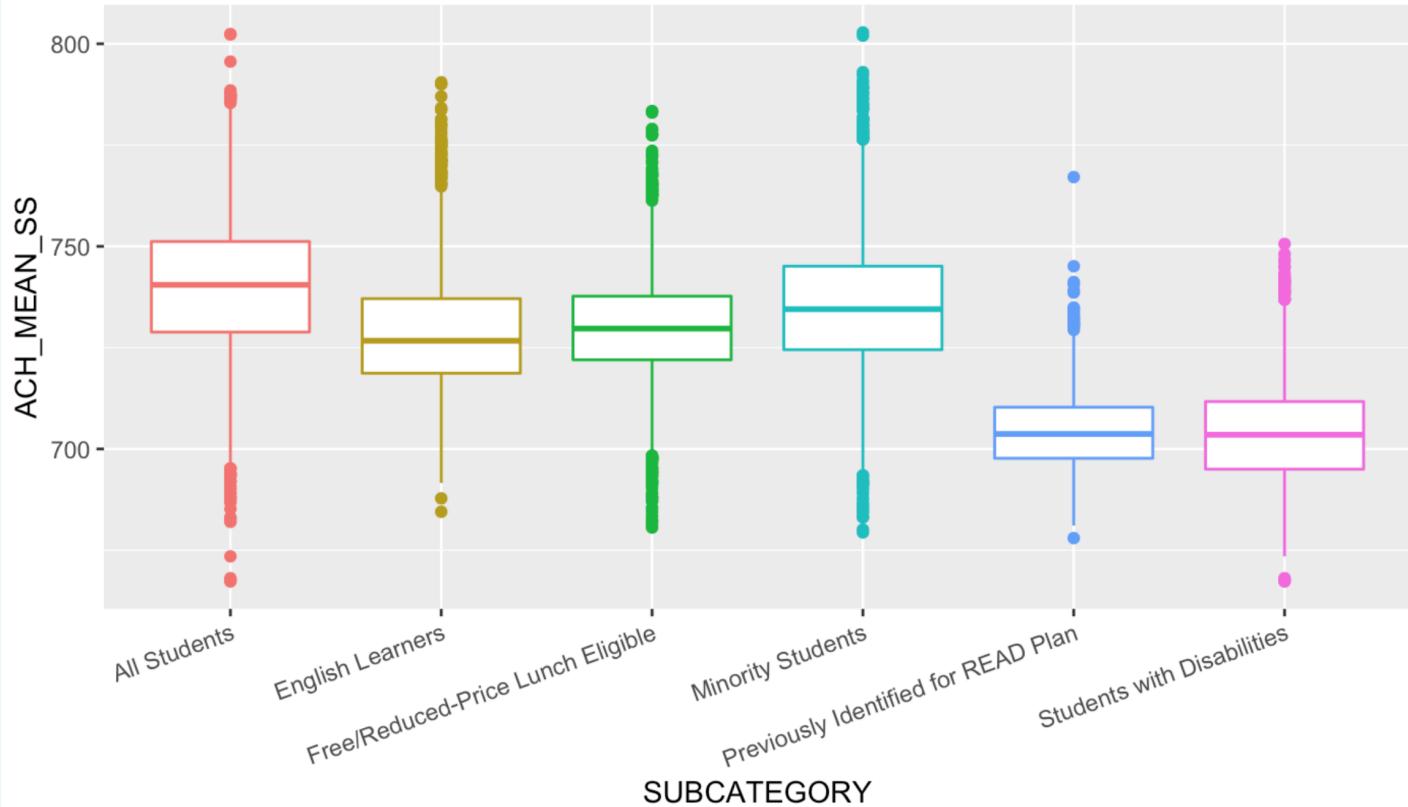
Results

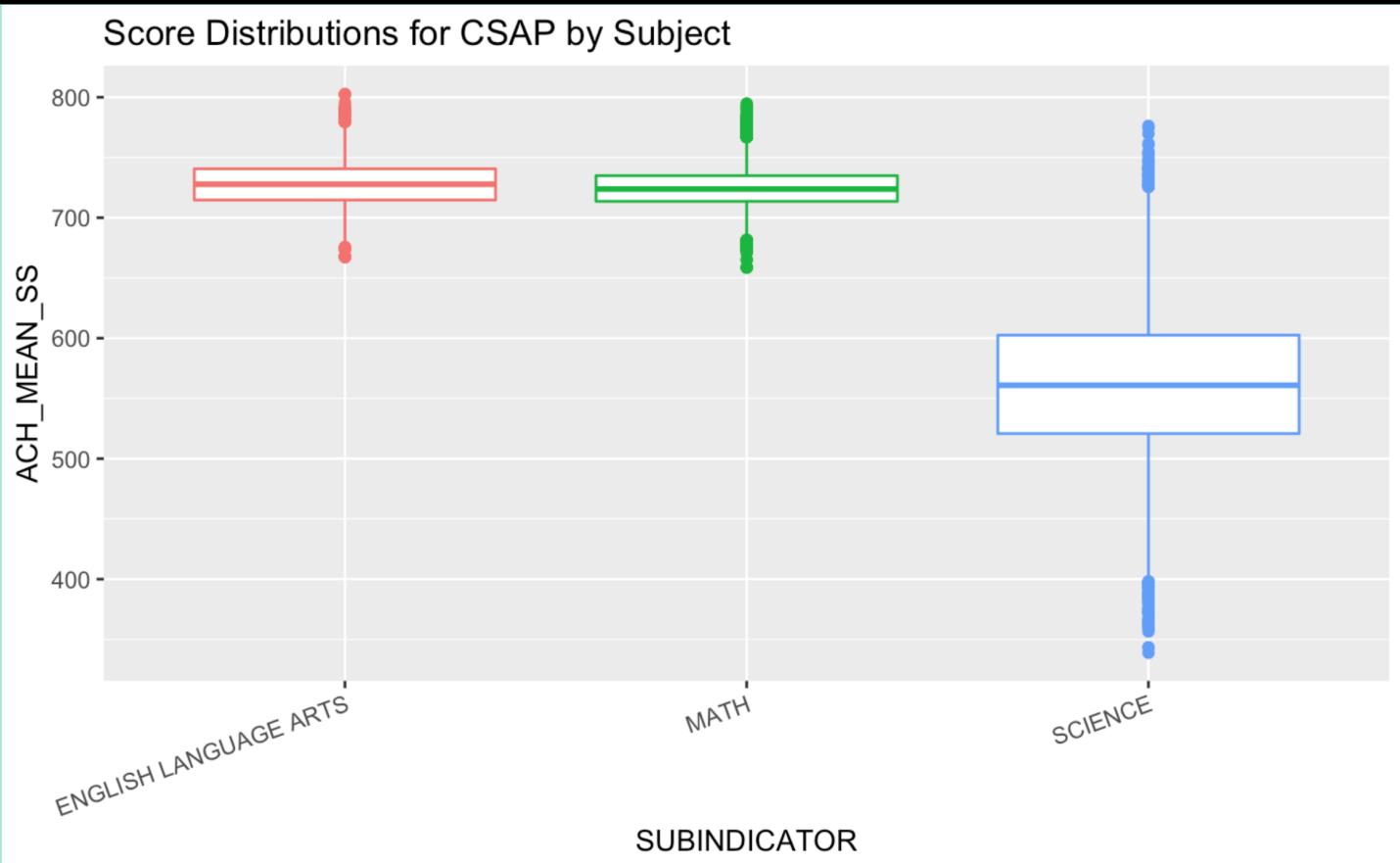
- Statistically significant differences between student subgroups
- Statistically significant differences between subjects

Student Subgroups

Subject

Score Distributions for CSAP English Language Arts Exams





Random Forest vs. K-Nearest Neighbors

Random Forest produced a lower error rate (RMSE) when compared to K-Nearest Neighbors in trying to predict a school's average CMAS exam scores.

Scores

Clusters

Predicting CMAS Scores Using RF

Science

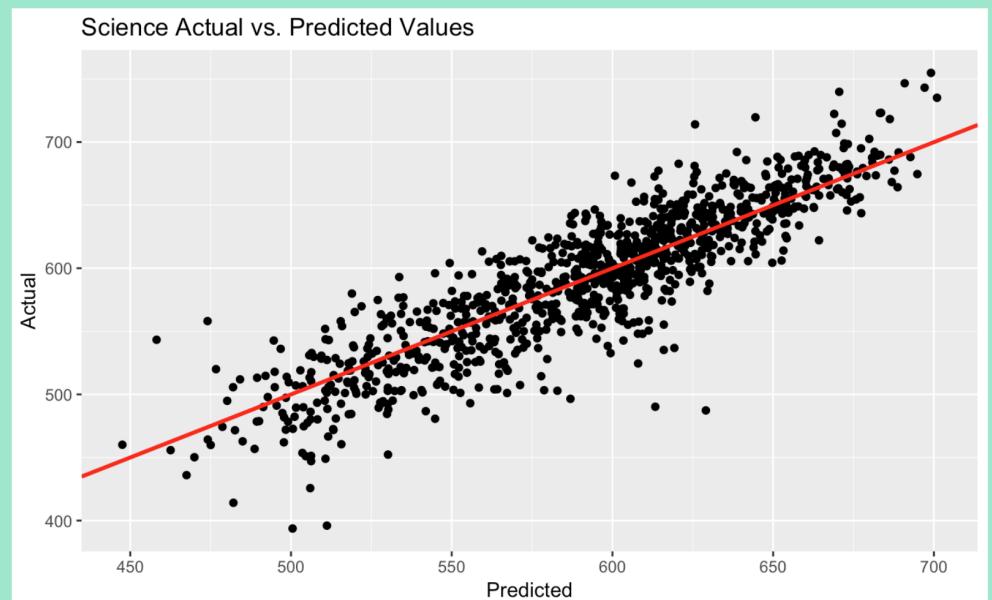
- RMSE: 28.1
- Range: 393-754

Math:

- RMSE: 7.6
- Range: 679-777

English:

- RMSE: 8.3
- Range: 667-795



Predicting Clusters

Accuracy: 97%

Kappa: 96.6%

The clusters are reproducible and can be predicted using machine learning.

Confusion Matrix and Statistics

		Reference		
		1	2	3
Prediction	1	798	13	11
	2	9	690	5
	3	6	7	731

Overall Statistics

Accuracy : 0.9775
95% CI : (0.9706, 0.9832)
No Information Rate : 0.3581
P-Value [Acc > NIR] : <2e-16

Kappa : 0.9662
McNemar's Test P-Value : 0.4697

Statistics by Class:

	Class: 1	Class: 2	Class: 3
Precision	0.9708	0.9801	0.9825
Recall	0.9815	0.9718	0.9786
F1	0.9761	0.9760	0.9805

Key Points

Are there differences in student groups?

YES

Are there differences by subject?

YES

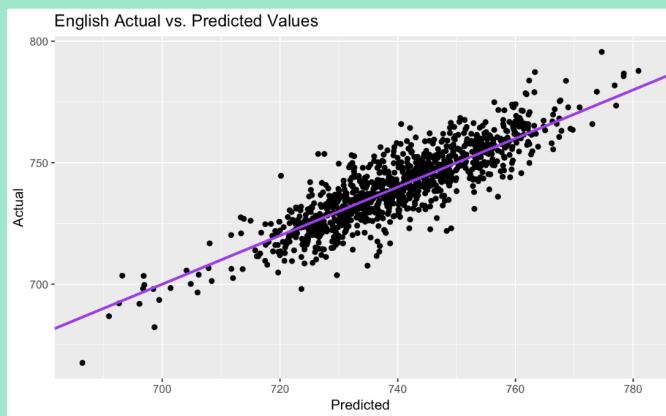
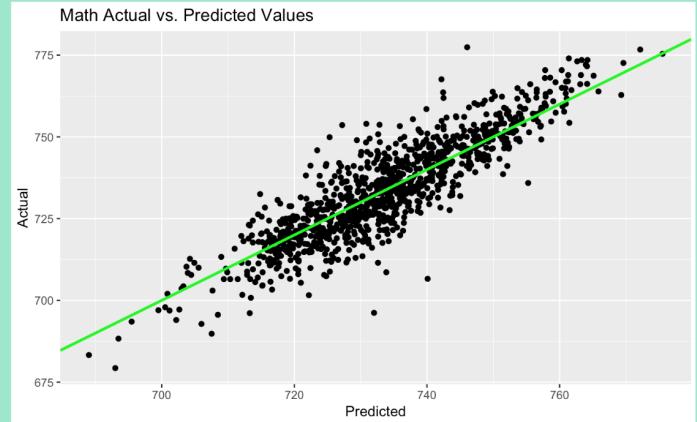
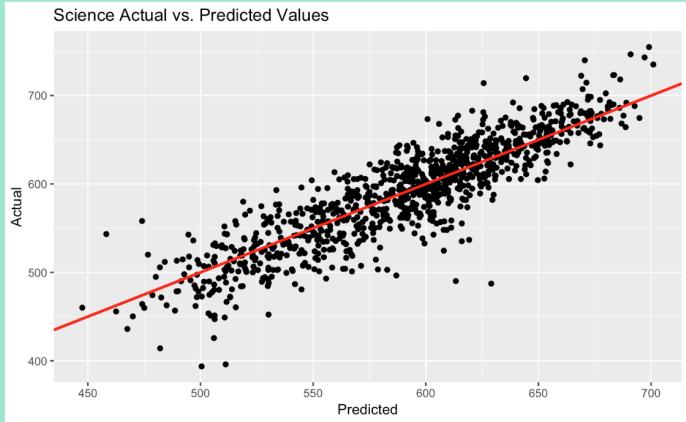
Can the scores be predicted?

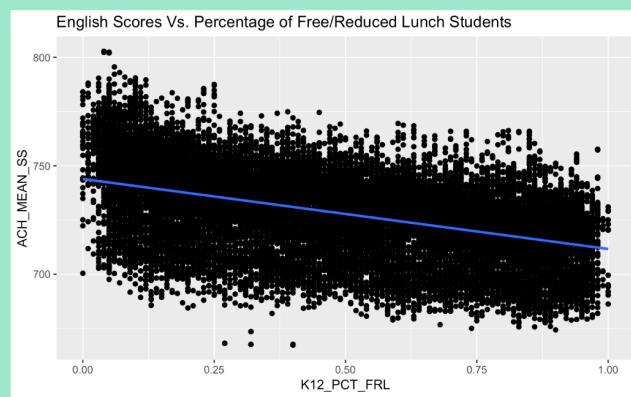
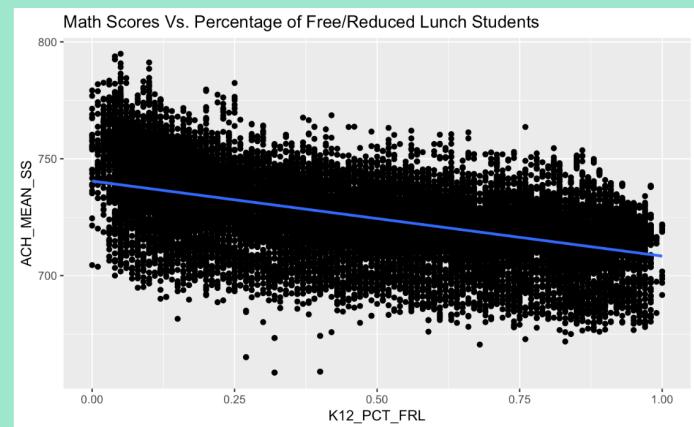
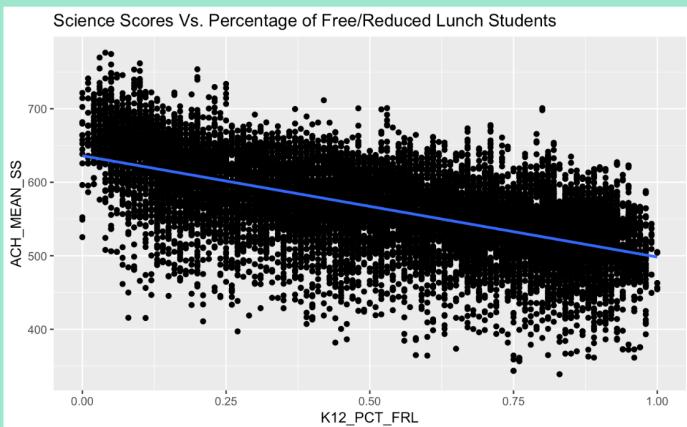
SOMEWHAT

Predicting

FRL

***Final
Thoughts***





*How can we help these students in ways
that go beyond what we do in the walls of
the school building?*