

BIOS 611 Final Project

Jeffrey Ayers

December 5, 2022

Education Trends for the United States

(Note: if you are reading the PDF of this file from the git repository, it may be out of date. It is included in the repo for convenience but the canonical writeup is in the org file in the repo. See the README instructions for how to build this pdf with make).

Introduction

As a Ph.D. student in the mathematics department instruction is a big part of my job, and one that I've become increasingly passionate about. Having grown up with a parent who is a teacher, and now becoming one I have developed an growing interest in seeing what factors lead to student success, especially from a monetary standpoint. I was fortunate to grow up in a state that places high value on education (Massachusetts) and have directly seen the benefits of a large instructional budget, but what about nationwide? Can we see any direct impact on student performance based on this? Moreover, what other hidden trends can we see from data? In this writeup I analyze data from the National Center for Education Statistics to see what trends occur in the US, and the Northeast.

Preliminary Figures

The data was retrieved from Kaggle: US Education Dataset, and minor cleaning was done to it; namely converting column names to lowercase, and removing the underling strings corresponding to state names. These data contained a wide data set of columns corresponding to states, years and a multitude of statistical information: Total enrollment, financial information, test scores from the NAEP exam for grades 4 and 8, as well as expenditure. For example, in the year 2011 we can see what the total enrollment was for each state.

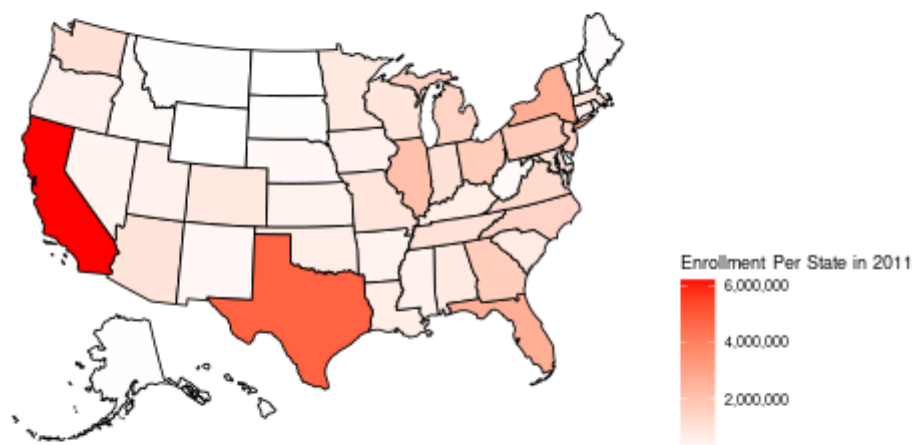


Figure 1: Total student enrollment per state in the year 2011

Correspondingly we can also view the total expenditure in dollars of each state, the total instructional expenditure, and the average NAEP math scores of each student in the states for grades 4 and 8:

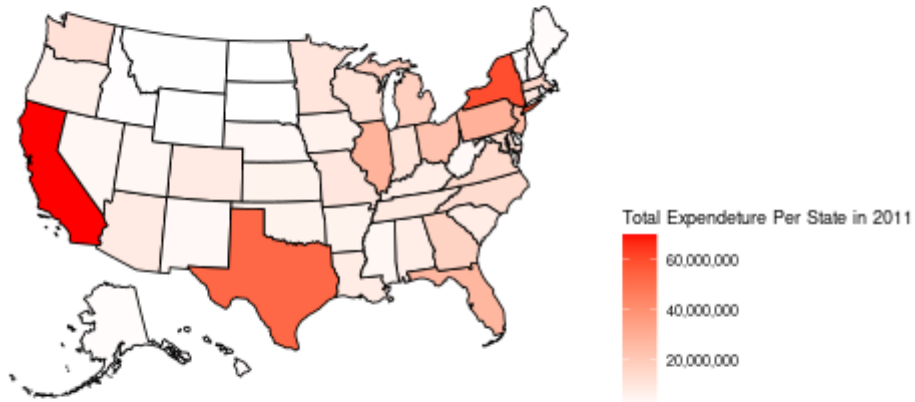


Figure 2: Total expenditure per state in the year 2011

It's no surprise that enrollment seems correlated with total expenditure: The more students the more money needs to be put in.

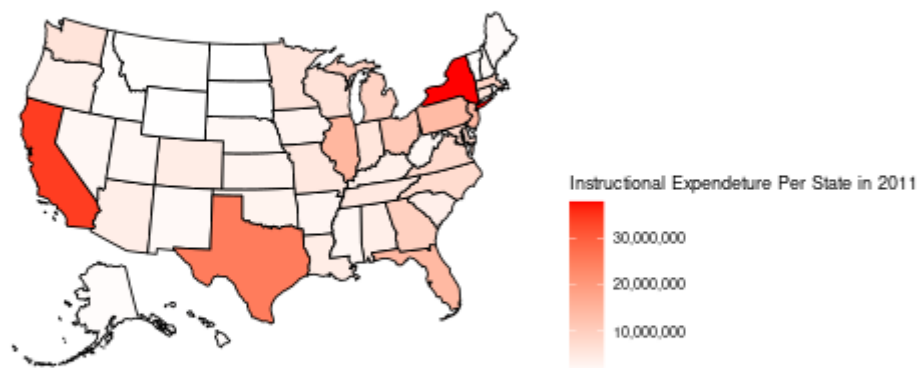


Figure 3: Total instructional expenditure per state in the year 2011

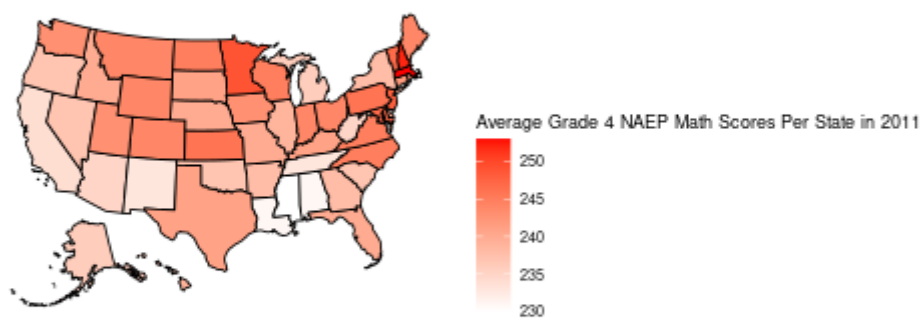


Figure 4: Average math score for grade 4 per state in the year 2011

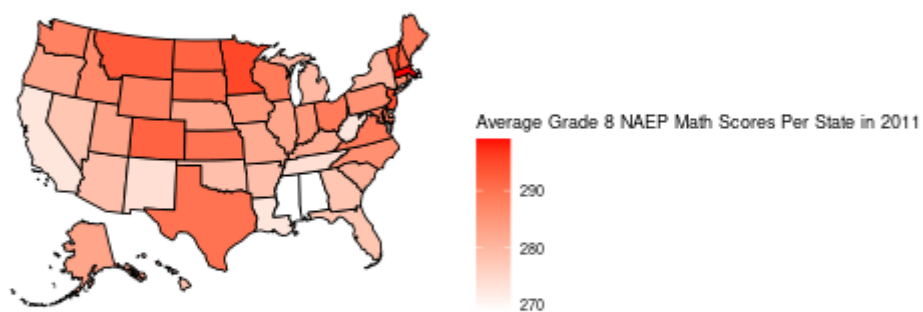


Figure 5: Average math score for grade 8 per state in the year 2011

There is some slightly better information. New York has the highest ratio at 14.13 dollars per student, for an

Analysis for New England

We can clearly see that Massachusetts and Connecticut far outpace the rest of the states.

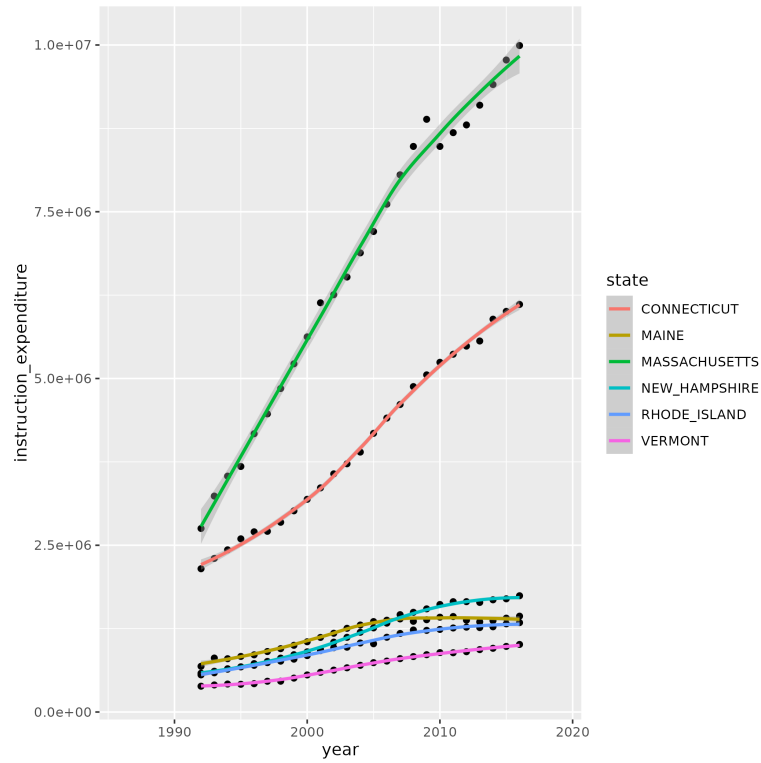


Figure 7: Instruction expenditure per year for New England

Looking at some more information we have a histogram of the average math scores for grade 4 students in New England:

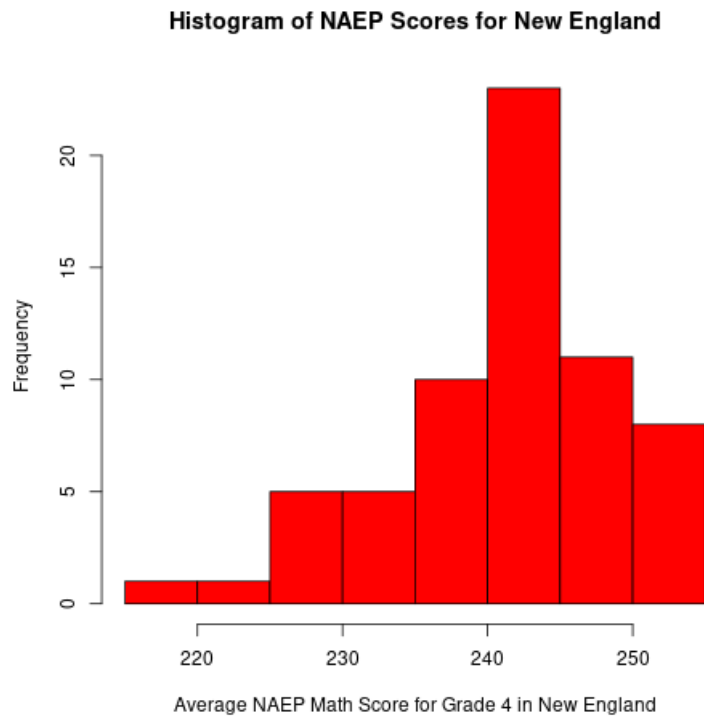


Figure 8: Histogram for New England scores in 2011

Now let's see how this compares with the rest of the US:

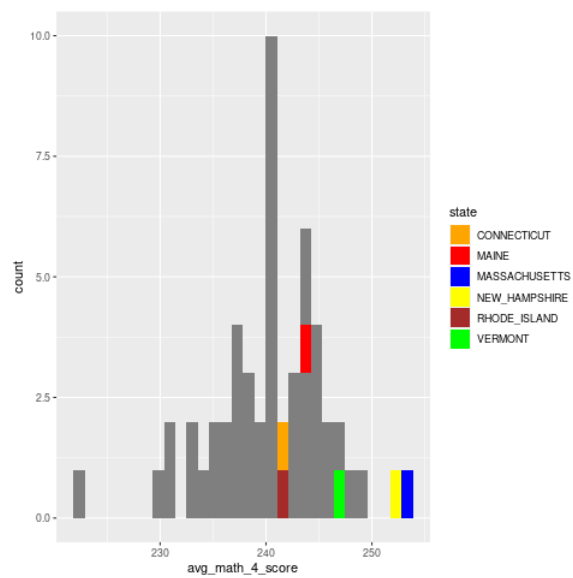


Figure 9: Histogram for US, with New England states colored in 2011

Here we get a good idea how much New England states out perform others, at least for 2011. We can perform a similar histogram with all years and each test to truly see some interesting data:

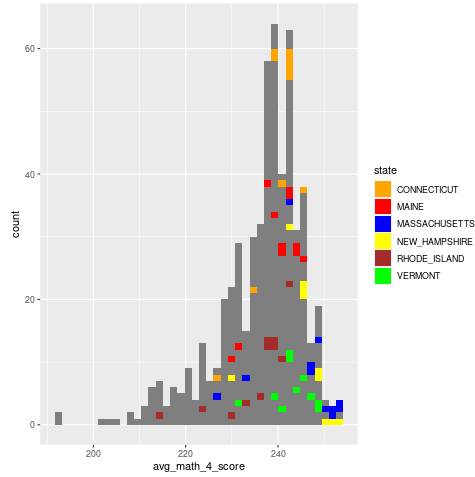


Figure 10: Histogram for US, with New England states colored for average grade 4 math score

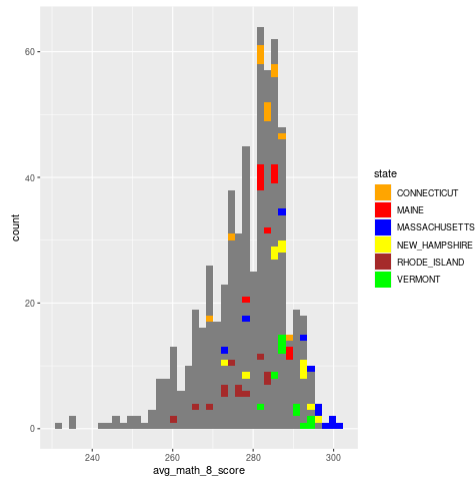


Figure 11: Histogram for US, with New England states colored for average grade 8 math score

While there is nothing definitive certainly we can see that apart from Rhode Island, New England students perform quite well on NAEP tests when they're offered.

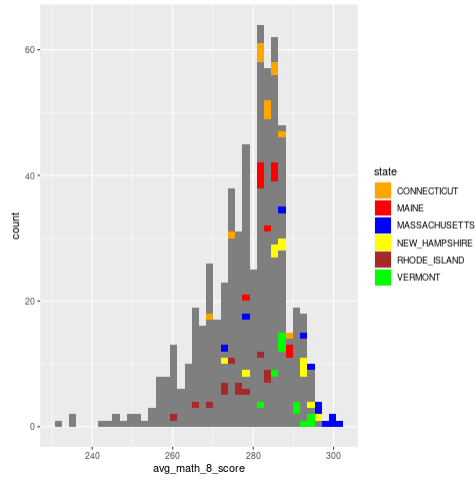


Figure 12: Histogram for US, with New England states colored for average grade 8 reading score

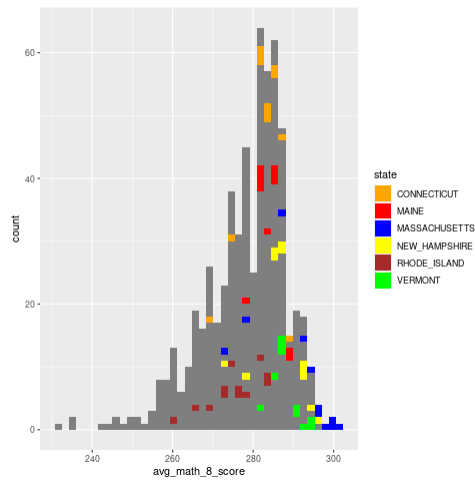


Figure 13: Histogram for US, with New England states colored for average grade 4 reading score

In an effort to see how much these states differ from each other, we perform a PCA. Given that these data are incomplete: various years are missing enrollment, test data, etc. we need to impute the data to ensure a valid PCA. My first attempt was to fill in all NA values with 0. The result is a quite bad association of states:

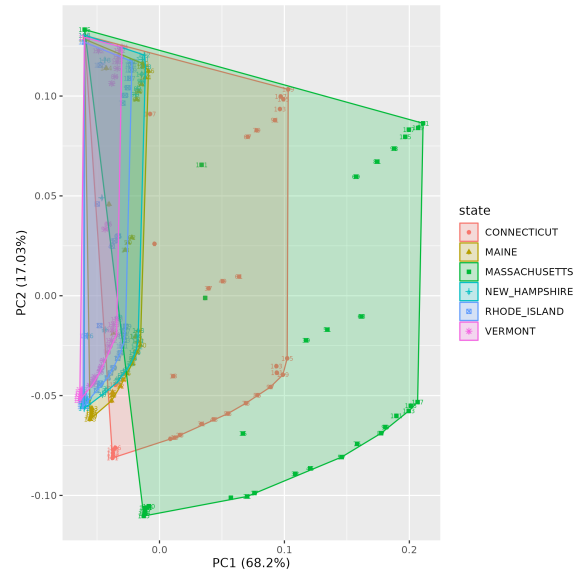


Figure 14: PC1 and PC2 for New England with 0's for NA values

However using the missMDA library, which imputes the missing values for us before we preform a PCA, we can see some nice association of states:

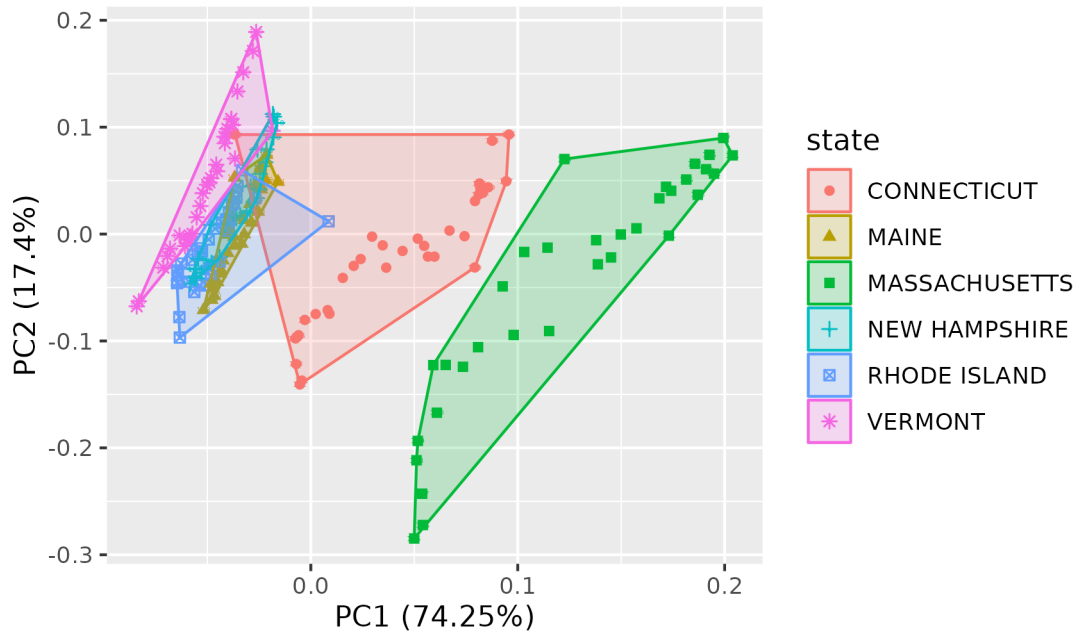


Figure 15: Imputed PC1 and PC2 for New England

With the imputed values we can see that not only do the proportional of variance increase, but we have more well defined clusters for Massachusetts and Connecticut, as well as Vermont. Moreover, the PCA tells us that the difference between Massachusetts is much more different than any other New England state.