

Project Writeup - Nelson Pollard

 [nwea/ai-academy-nelson-pollard](https://github.com/nwea/ai-academy-nelson-pollard)

Description

Use machine learning to predict distribution of students across Achievement Level Descriptors (ALDs). Explore what features increase or decrease classification accuracy, with an interest in how political affiliation of a state may affect ALDs. Note that to simplify data collection and processing, only states that have adopted [Common Core](#) were included.

The ALDs are:

1. Below Basic (BB)
2. At Basic (BA)
3. At Proficient (PR)
4. At Advanced (AD)

Hypothesis

The original hypothesis was that political affiliation of a state, as well as the number of years since the last standard setting, would greatly influence the distribution of students across ALDs for a given year/grade/subject. The idea was that there was some sort of bimodal “ideal” distribution, with one normal curve for right-leaning states and another for left-leaning states. There might be some drift between standard settings, but immediately after a standard setting was performed we’d see the distribution across ALDs “snap back” into an approximation of the “ideal” distribution.

Features

- State
- Year
- Political affiliation (Republican, Democrat, Independent, or Unknown) of candidate voted for by state in most recent presidential election
- Political affiliation (Republican, Democrat, Independent, or Unknown) of governor of state
- Years since most recent standard setting

Data Sources

- ALDs: [NAEP Data Service API](#)
- Presidents: [GitHub - zonination/election-history: US Presidential Elections since 1789](#)
- Governors: [BP Governor \(state executive office\) - Ballotpedia](#)

Algorithms

- Linear regression: [sklearn.linear_model.LinearRegression](#)
- Multi-layer Perceptron regressor: [sklearn.neural_network.MLPRegressor](#)
 - hidden_layer_sizes=(30,30,30)
 - max_iter=1000

Results

Machine learning

Linear regression

Linear regression resulted in accuracy rates just below 80% across the board, regardless of the combination of features in/excluded.

Truncated output:

```
1 -----
2 LinearRegression()
3 POLITICS | STATES | YEARS | YEARS SINCE STANDARD SETTING
4 0.7877199030960381
5 -----
6 LinearRegression()
7 POLITICS | STATES | YEARS | !YEARS SINCE STANDARD SETTING
8 0.7854942455283795
9 -----
10 LinearRegression()
11 POLITICS | STATES | !YEARS | YEARS SINCE STANDARD SETTING
12 0.7843786682897814
13 -----
14 LinearRegression()
15 POLITICS | STATES | !YEARS | !YEARS SINCE STANDARD SETTING
16 0.7922950734581131
17 -----
18 LinearRegression()
19 POLITICS | STATES | YEARS | YEARS SINCE STANDARD SETTING
20 0.7859662929352996
21
22 ...
```

Multi-layer Perceptron

The highest accuracy was achieved when including the “State” feature and excluding the “Year” feature, just under 97%, with the inclusion of the number of years since the last standard setting increasing accuracy by less than half a percent:

```
1 -----
2 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
3 !POLITICS | STATES | !YEARS | !YEARS SINCE STANDARD SETTING
4 0.9676490507049749
5 -----
6 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
7 POLITICS | STATES | !YEARS | !YEARS SINCE STANDARD SETTING
8 0.9677645069243545
9 -----
10 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
11 !POLITICS | STATES | !YEARS | YEARS SINCE STANDARD SETTING
12 0.9691524415578044
13 -----
14 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
15 POLITICS | STATES | !YEARS | YEARS SINCE STANDARD SETTING
16 0.9693550880790959
```

Excluding the “State” feature yielded accuracy in the high 80s:

```
1 -----
2 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
3 !POLITICS | !STATES | !YEARS | !YEARS SINCE STANDARD SETTING
```

```

4 0.8711898454740707
5 -----
6 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
7 POLITICS | !STATES | !YEARS | !YEARS SINCE STANDARD SETTING
8 0.8822530246594544
9 -----
10 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
11 !POLITICS | !STATES | !YEARS | YEARS SINCE STANDARD SETTING
12 0.8720625957121539
13 -----
14 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
15 POLITICS | !STATES | !YEARS | YEARS SINCE STANDARD SETTING
16 0.8783474478269986

```

Including “Year” decreased accuracy to the 60s-70s range, likely due to overfitting:

```

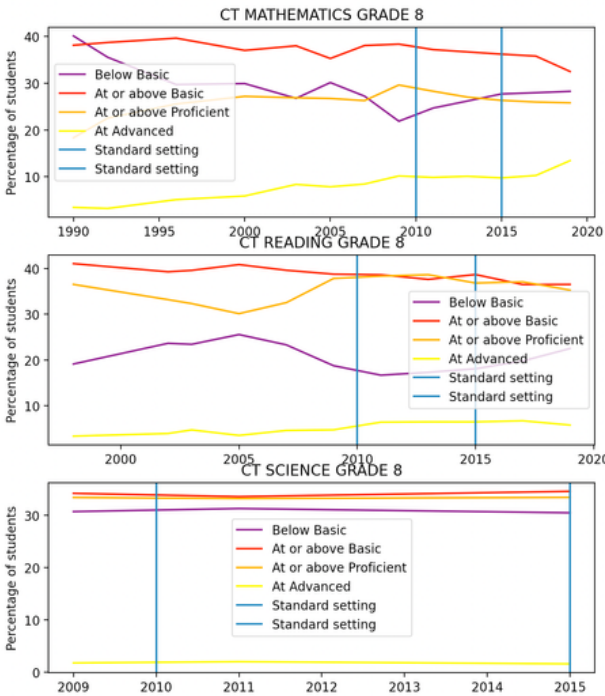
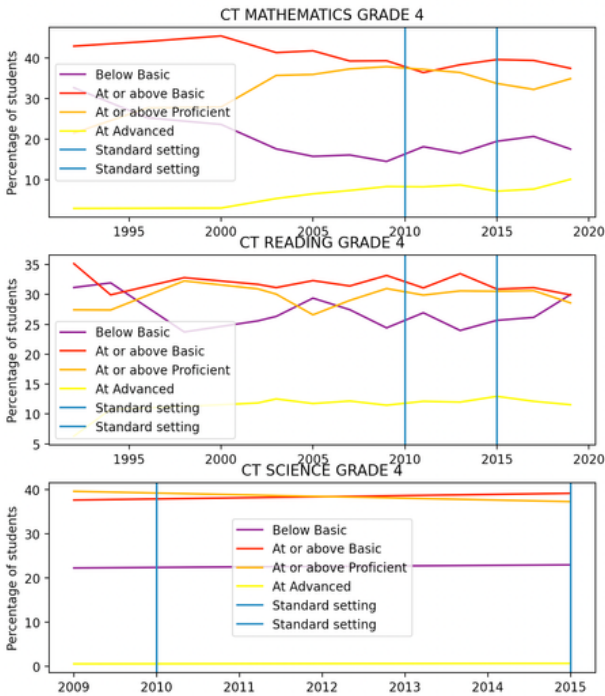
1 -----
2 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
3 !POLITICS | STATES | YEARS | !YEARS SINCE STANDARD SETTING
4 0.6893920744030189
5 -----
6 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
7 POLITICS | STATES | YEARS | !YEARS SINCE STANDARD SETTING
8 0.6590560372043667
9 -----
10 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
11 !POLITICS | STATES | YEARS | YEARS SINCE STANDARD SETTING
12 0.6123630826575253
13 -----
14 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
15 POLITICS | STATES | YEARS | YEARS SINCE STANDARD SETTING
16 0.7103966368277528
17 -----
18 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
19 !POLITICS | !STATES | YEARS | !YEARS SINCE STANDARD SETTING
20 0.7718847283209895
21 -----
22 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
23 POLITICS | !STATES | YEARS | !YEARS SINCE STANDARD SETTING
24 0.6974333051365792
25 -----
26 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
27 !POLITICS | !STATES | YEARS | YEARS SINCE STANDARD SETTING
28 0.6973731449238002
29 -----
30 MLPRegressor(hidden_layer_sizes=(30, 30, 30), max_iter=1000)
31 POLITICS | !STATES | YEARS | YEARS SINCE STANDARD SETTING
32 0.6813192856326488

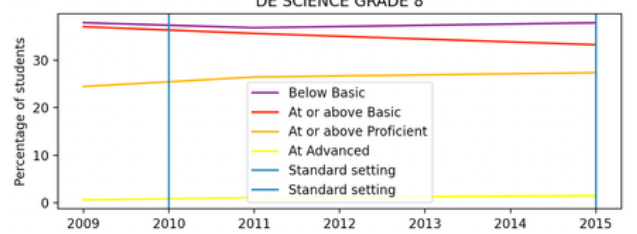
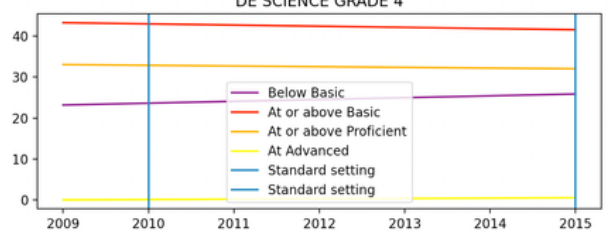
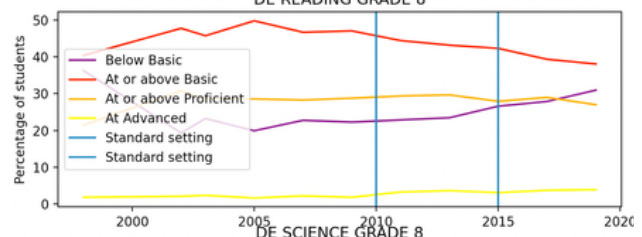
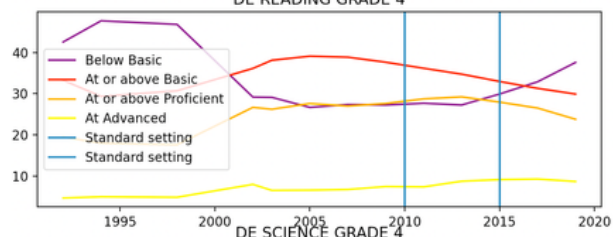
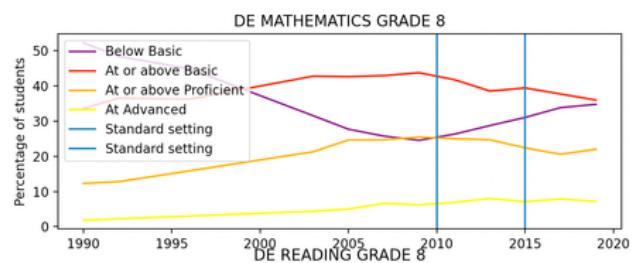
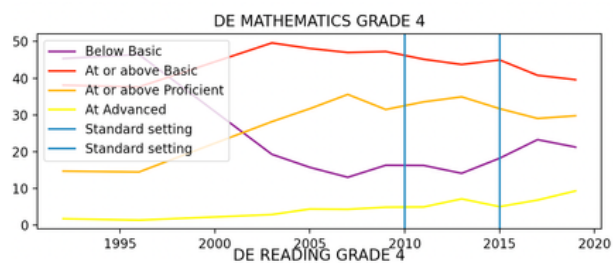
```

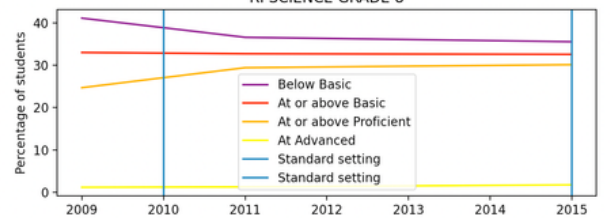
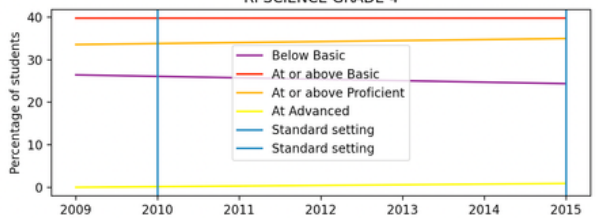
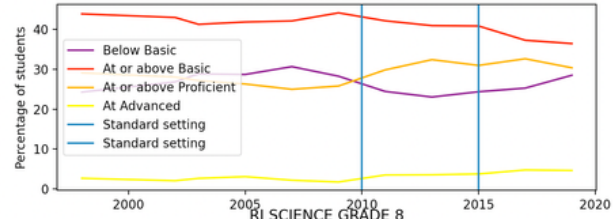
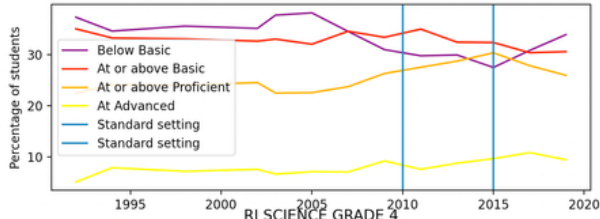
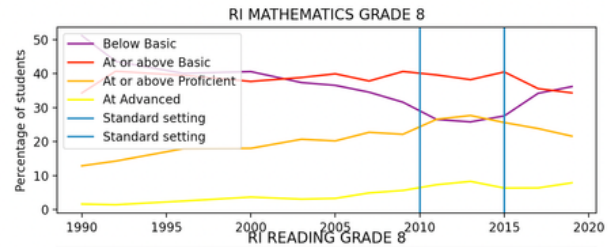
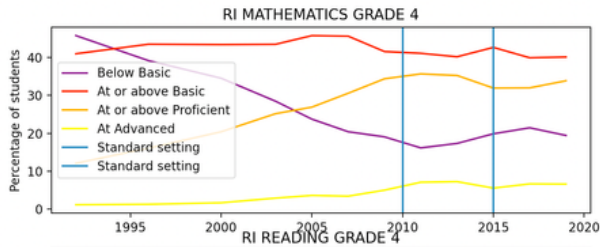
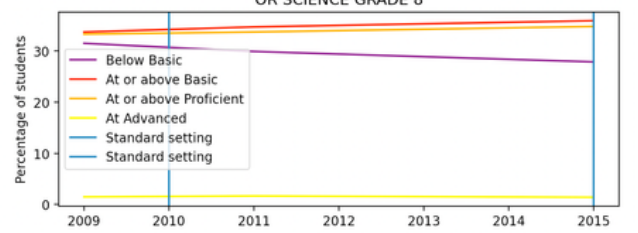
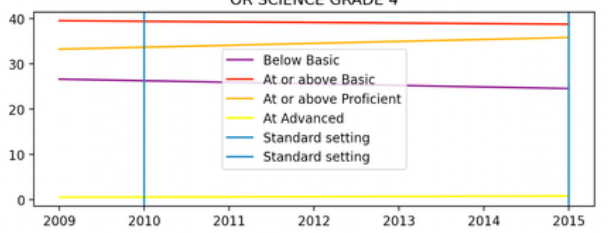
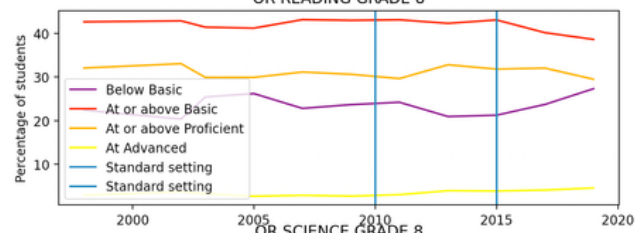
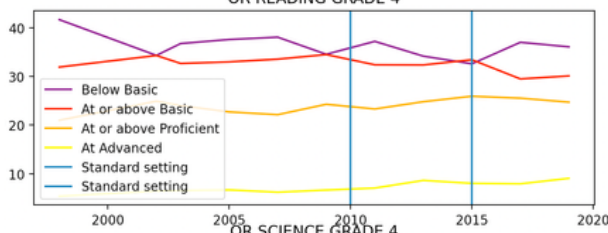
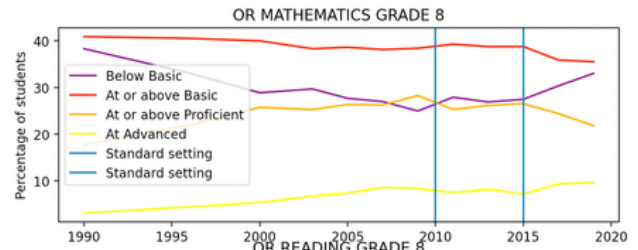
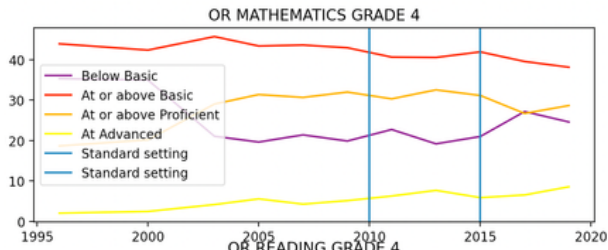
The inclusion or exclusion of political affiliation and the years since last standard setting appeared to have a negligible affect on accuracy.

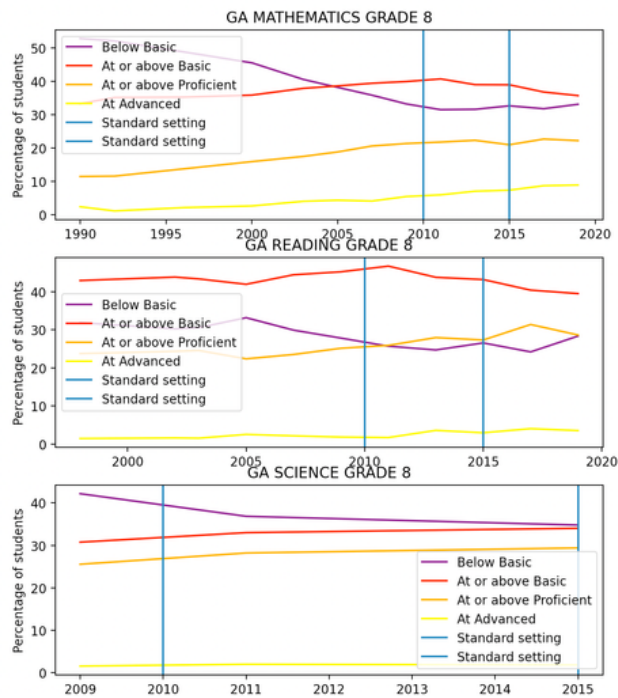
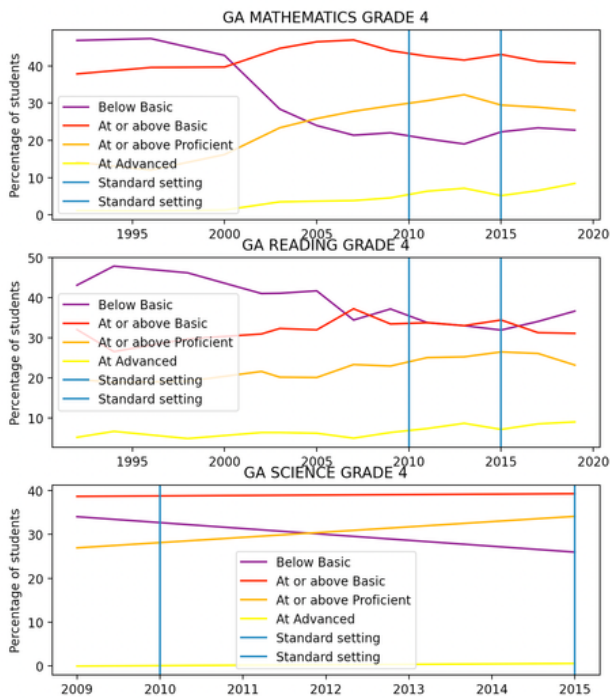
Pretty graphs

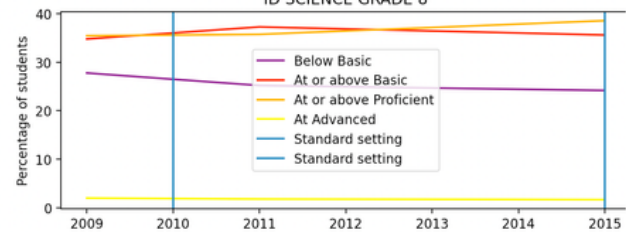
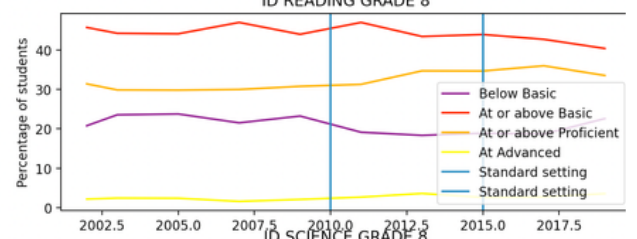
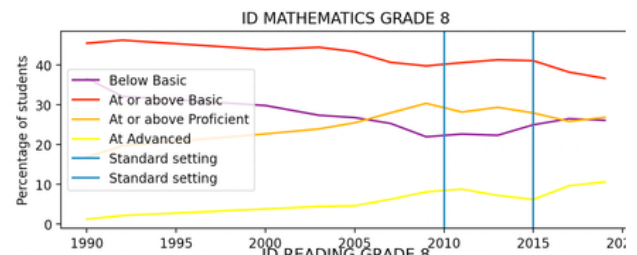
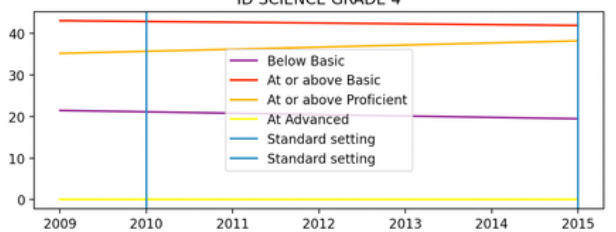
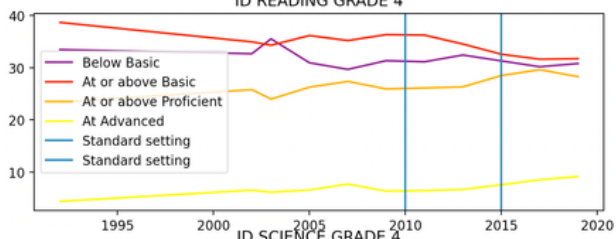
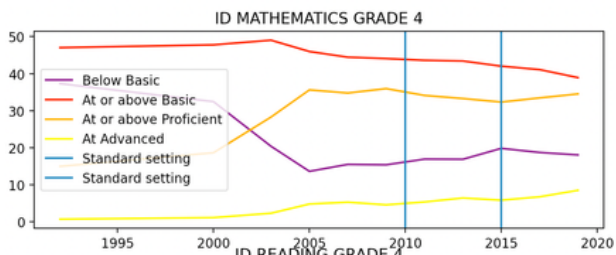
In addition running machine learning algorithms against the data, graphs were generated to visually check for trends, especially in relation to the years since last standard setting:

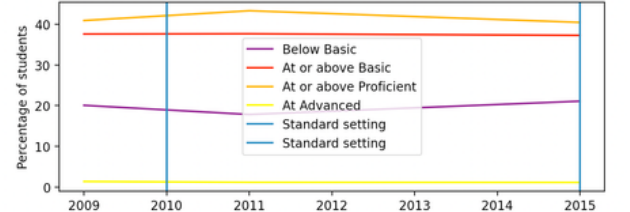
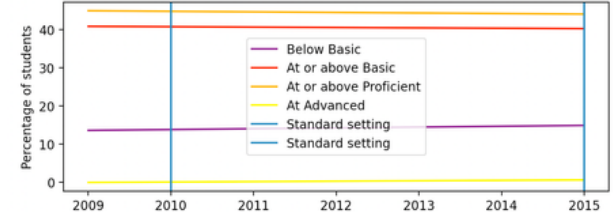
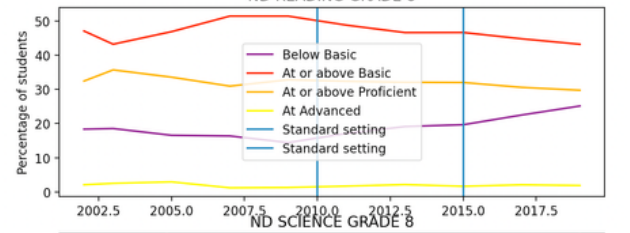
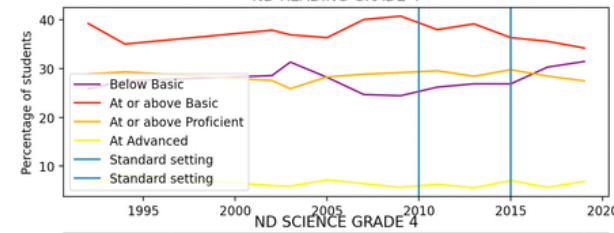
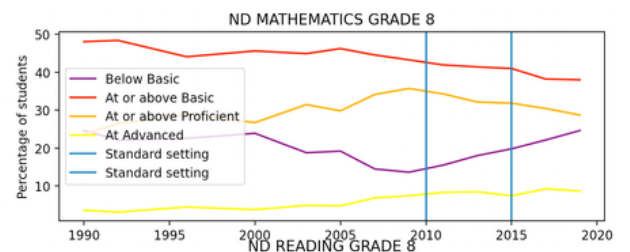
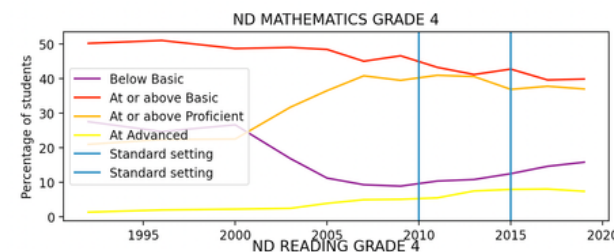
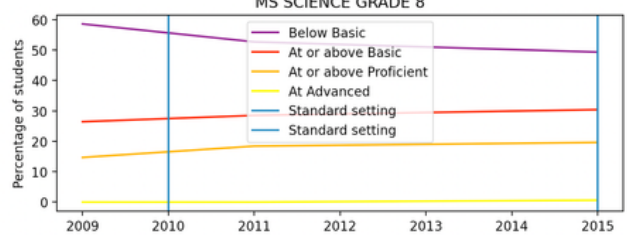
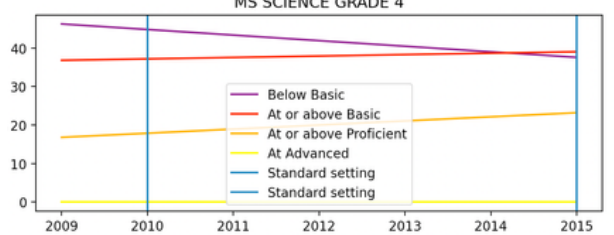
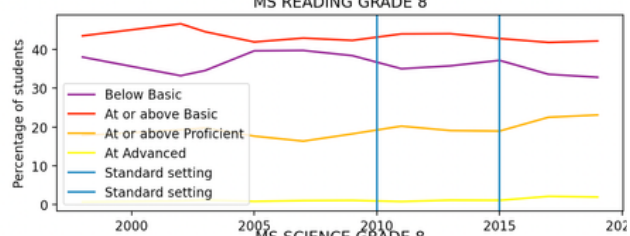
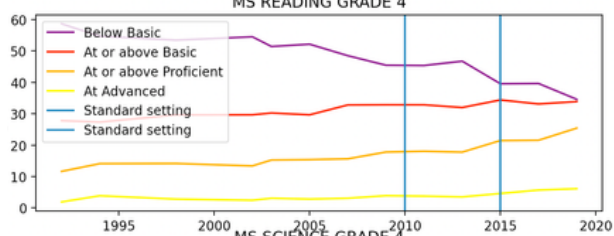
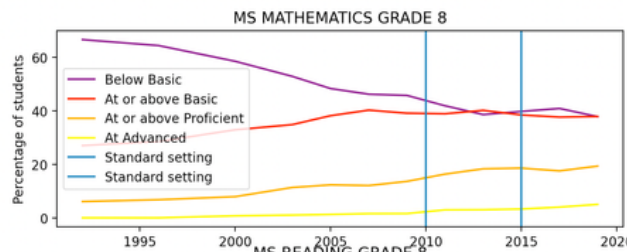
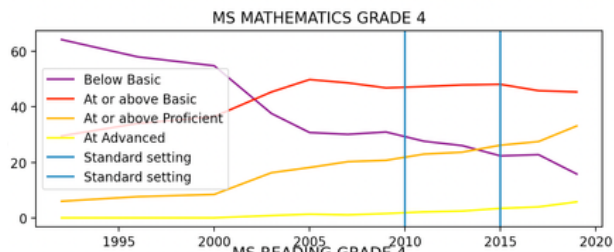


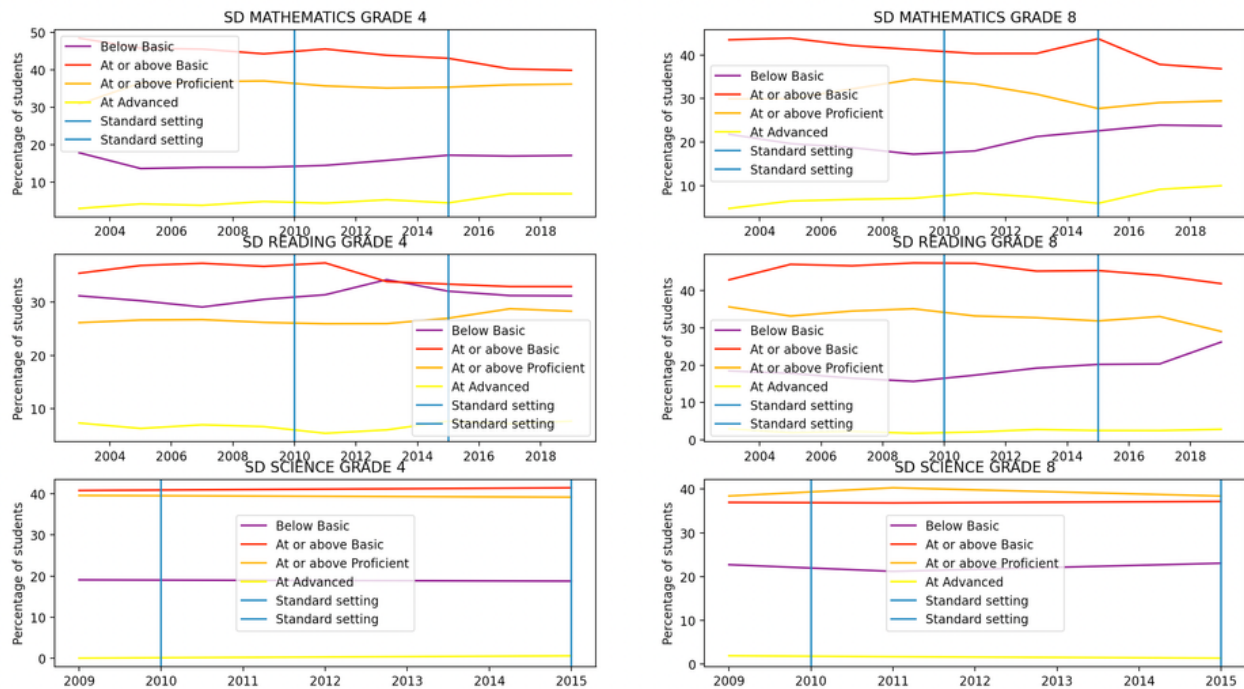




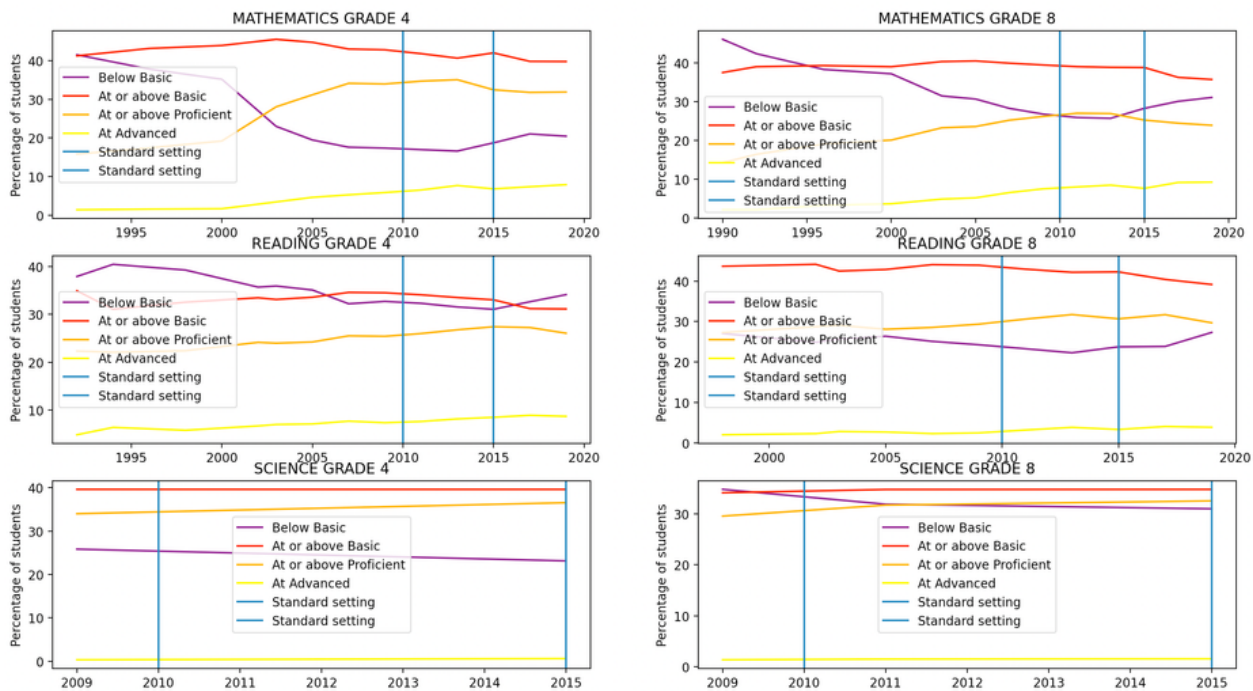




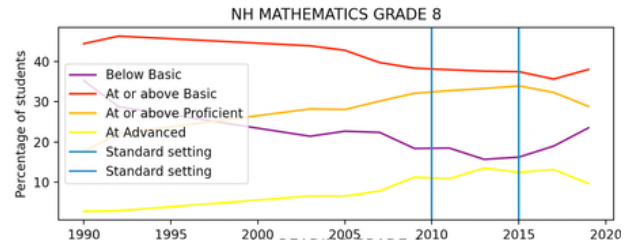
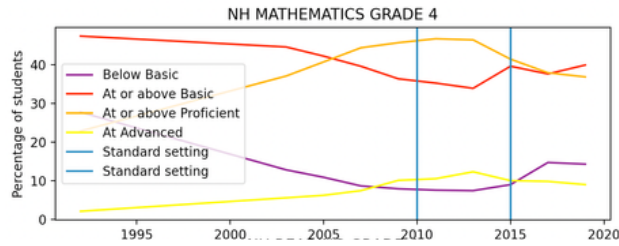
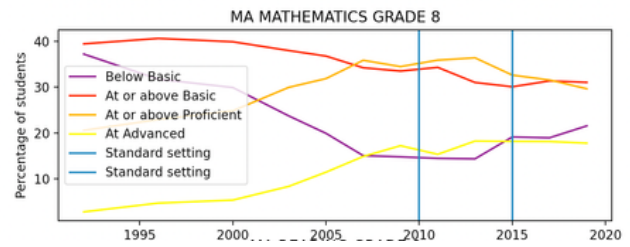
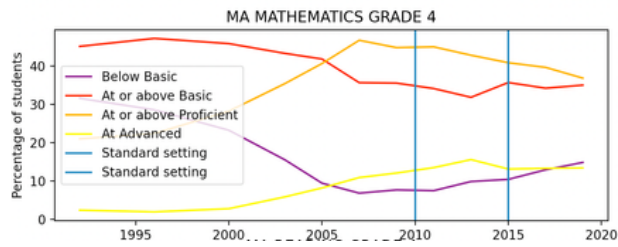




Since Common Core states have all adopted the same standards, an average across all states was also plotted:



Two observable patterns are the relative sparsity of "Science" data and the almost universally low percentage of students in the "At Advanced" ALD. The only exceptions to this latter pattern are:



Conclusion

The graphs did not show any obvious “snapping back” to an “ideal” distribution following standard settings. There is also not an obvious difference in the distributions of the most left-leaning states versus the most right-leaning states. The fact that these are Common Core states, and thus are sharing standards, might explain why we do not see any obvious skewing in relation to political affiliation. Perhaps states that set their own standards display such skewing. Of course, it’s also possible that processing the data differently, using different machine learning algorithms, and/or visually plotting in a different manner would bring forth different results. But for this given approach to this set of data, it must be concluded that we have not proven the hypothesis.