# Change in State: using Markov chains to explore national economic mobility and the evolving American Dream

Daniel Ayasse, Emily Myers, Christian Schmidt, Allison Schwam

July 27, 2016

**Abstract**

By classifying the economic status of children and their parents into equal quintiles, it is possible to use Markov chains to provide general statements about intergenerational economic mobility. Using income tax records from over 40 million children and their parents, we examine the evolution of the probability of achieving the "American Dream" over three generations. We then explore areas found to have greater upward mobility and attempt to identify contributing factors and trends. Our findings show a moderately high, negative correlation between attainability of the American Dream and race.

## 1 Introduction

The United States is known historically as the land of opportunity, where anyone can change their fortunes regardless of their economic status. However, we show that this description is not always merited. Intergenerational mobility varies drastically by area, and while some children are able to climb the rungs of the economic ladder, others simply aren't given the opportunity.

## 1.1 Aims

In this paper we analyze economic mobility in each state using data taken from federal income tax records and separated first by state and then by Commuting Zone (CZ). CZs are defined by "geographical aggregations of counties that are similar to metro areas but cover the entire U.S., including rural areas."[2] We define economic mobility as the movement of people from one economic class to another, with each 20th percentile as a different economic class. We analyze the movement of children from their parents' economic class to the economic class they have reached at adulthood, age 30. We are specifically interested in the movement of children from the lowest economic class to the highest. We call this movement the achievement of the "American Dream". When we refer to the American Dream, we mean $P(child\,Q_5|parent\,Q_1)$. This reads as the probability that a child ends up in the 5th quintile, given that the parent is in the 1st quintile.

By separating the data by state, we are able to analyze which areas of the United States are the best and worst places to realize the American Dream. We also attempt to identify any factors that lead to discrepancies between different states' American Dream attainability. Factors we analyze include income disparity, race, population density, education, and religious affiliation.

## 1.2 Past work

### 1.2.1 The PEW Trust EMP: Urahn et al.

In the Economic Mobility Project (EMP), Urahn et al. examined generalized mobility for the entire United States. By separating economic status into quintiles, they created a series of matrices which they then analyzed in order to determine the economic enviroment has evolved over one generation.

In their work, they found that mobility is affected by a number of different factors. Noteable observations relevant to our work include [3]:

- Most Americans are able to move up the economic ladder and surpass both their parents in family wealth and income.

- Race is a factor in economic mobility. It is more difficult for African Americans to surpass their parents' family wealth and income than for whites to surpass their parents' family wealth and income. It is also more likely for African Americans to fall in the economic ladder or stay stuck on the bottom rung.

- College education is a factor in economic mobility. It was shown that a 4-year college degree both prevents a person's fall down the economic ladder as well as promotes the person's movement up the ranks.

### 1.2.2 NBER: Chetty et al.

In this paper, funded by the National Bureau of Economic Research (NBER), Chetty et al. used federal tax data from one generation to approximate the

linear relationship between parent and child economic status.

They found that intergenerational mobility varies substantially across different states (and CZs). In their exploration of factors that potentially affect mobility, they found that areas with high upward mobility are defined by [2]:

- Less residential segregation

- Less income inequality

- Better elementary schools

- Better social capital, which is to say that economic transactions are characterized by trust and reciprocity

- Greater family stability

### 1.2.3  Extending past work

The PEW Trust's EMP presents valuable results, but their findings are for the entire United States and are therefore less precise. For example, a family in the 5th quintile in a state might have top economic status relative to the rest of their state, but not on the national level.

Meanwhile, the NBER paper differentiates between extremely specific areas but does not extend to predict mobility for any future generations. They also specify several measures of absolute mobility, but the one that interests us most is the probability of a child being born in the bottom quintile ending up with income in the top quintile: what we refer to as the American Dream. We are particularly interested in the evolution of this specific type of mobility over generations.

In our project we combine aspects of both papers, using the specificity of the NBER project and the quintile system of the EMP: applying Markov chains to the NBER tax data in order to predict future mobility for each state. We do this by compiling CZ for each state and creating a state-defined mobility matrix.

## 2  Data

In the course of our work it has become necessary to make a certain set of definitions and assumptions, almost all of which were defined by either of the previous works.

We begin by defining intergenerational mobility. Our core data from federal tax records between 1996 and 2012, compiled by Chetty et al [1]. Family income is measured between 1996 and 2000, when the children in question are in their teens, and the child's income is measured in 2011-2012, when they are roughly 30 years of age and presumably established in life. This is the same data used by the NBER. The parents are defined as the first person or persons to claim the child on a tax form. This data assumes that the child

- is in possession of a valid Social Security or Taxpayer Identification Number

- was born between 1980 and 1991

- is a United States citizen as of the year 2013

## 2.1 Limitations

When testing the data, we noticed that the rows of our eventual probability matrices would not sum to 1, which is a significant issue when desiring a probability space. For example, for Anchorage, Alaska we have the following data

| CZ Name | $P(Child\,Q1|$ $Parent\,Q1)$ | $P(Child\,Q2|$ $Parent\,Q1)$ | $P(Child\,Q3|$ $Parent\,Q1)$ | $P(Child\,Q4|$ $Parent\,Q1)$ | $P(Child\,Q5|$ $Parent\,Q1)$ |
|---------|-----------|-----------|-----------|-----------|-----------|
| Anchorage | 0.295 | 0.215 | 0.189 | 0.167 | 0.134 |

Table 1: Data from Anchorage, AK

We see that

$$0.295 + 0.215 + 0.189 + 0.167 + 0.134 = 1$$

However, the values that we see are rounded to three decimal places automatically by Excel. In truth, they are much longer decimals, anywhere from 6 to 15 decimal places long. If we take our sum out to 8 decimal places, we find that the sum is actually 0.99999998. To deal with this, we only use the values for calculation up to the three decimal places that are reported on the Excel sheet. This ensures that our rows all sum to one, as desired.

## 2.2 Measuring Intergenerational Mobility [2]

While there are some issues regarding bias in measures of intergenerational mobility, we use the income definitions of Chetty et al., who have shown that the age at which parental income is measured is irrelevant between the ages of 30 and 55, estimates stabilize as the child approaches their late twenties, and that for all practical purposes, using multiple years of data does not improve the estimate.

Again as with Chetty et al., we use a rank-rank system to compensate for any kind of bias. In this system children are compared to other children based on their incomes. Parents are also compared to each other on the basis of income, and they find that there is a linear relationship between a child's income rank and their parents' income rank, with a slope of 0.341. In other words, an increase of 10 percent in a parent's rank would indicate a 3.41 percent increase in their child's rank. In fact, as we will discuss later, many other factors are linearly related to parental income.

The linear relationship between child and parent ranks retains its linearity even within commuting zones, so by using the slope for any commuting zone

(CZ) in combination with an intercept that represents the expected rank for children in the bottom income group for that CZ, we can easily calculate the expected rank for any child also in that CZ. This value is referred to as absolute mobility at percentile p, where p is the percentile of the parental income distribution [2].

We can see in the data that both relative and absolute mobility vary highly by geographic location. Relative mobility is highest in the rural midwest, but absolute mobility is highest in Salt Lake City with $p = 46.2$.

Though Chetty et al. define several different measures of absolute mobility, as we have already discussed, the one we will be concentrating on is the probability that a child will be able to to rise from the bottom to the top quintile: the American Dream.

## 2.3   Omitting the D.C. Commuting Zone

We feel the need to omit the D.C. CZ from our analysis for a couple of reasons. Our main concern involves the size of the D.C. CZ, which comprises a large section of both Maryland and Virginia. In fact, the Maryland and Virginia populations in the CZ are both individually higher than the D.C. population in the CZ. By labeling this CZ "DC", we would force ourselves to create a new data point, one that notably ends up acting as an outlier with regard to our race analysis. While we cannot be sure, this behavior is almost surely due to the Maryland and/or Virginia populations, and so including D.C. as a CZ would offer a skewed perception of the area, as well as affect our analysis. For these reasons, we feel it is necessary to omit D.C. from our paper.

It is also important to note that the states of Maryland and Virginia have American Dream values that are affected by this grouping of certain areas with D.C. However, without more detailed data, we are unable to separate the CZ into its components, which would be the optimal solution.

# 3   Methodology

We started out with data for each CZ in America. Specifically, for each CZ zone, we were given

$$P_{CZ}(child\,ends\,up\,in\,quintile\,j\,|\,parent\,is\,in\,quintile\,i)$$

for $i, j \in [1, 5]$. However, we wanted data for each of the 50 states in order to more easily create maps and thus compare data across regions of the US. We used Maple to do this. Our first task was to compute the weight of each CZ with its respective state.

$$CZ_{weight} = \frac{CZ\ Cohort\ Population}{State\ Cohort\ Population}$$

And to get the state-wide probabilities, we compute a weighted sum,

$$P_{State}(child\,Q_j|parent\,Q_i) = \sum_{CZ_n \in State} CZ_{n,weight} \cdot P_{CZ_n}(child\,Q_j|parent\,Q_i)$$

where $n$ allows us to index through the CZs in the given state. We'll look at Delaware as a short example, since Delaware only has two CZs.

Let's compute $P_{Delaware}(child\,Q_3|parent\,Q_1)$. Delaware has two CZs: Wilmington and Dover. We'll deal with Wilmington first. Since Wilmington has a higher cohort population than Dover, we want and expect its weight to be higher.

$$Wilmington_{weight} = \frac{42343}{74367}$$
$$= 0.569$$

$$Dover_{weight} = \frac{32024}{74367}$$
$$= 0.431$$

Now we can use these weights to compute our desired probability.

$$P_{Delaware}(child\,Q_3|parent\,Q_1) = 0.569 \cdot 0.181 + 0.431 \cdot 0.186$$
$$= 0.183$$

Once we computed each $P_{State}(child\,Q_j|parent\,Q_i)$, for any given state, we then created a matrix of all the values. Let's use our Delaware example.

$$
\text{Parent}
\begin{array}{c}
\text{Child} \\
\begin{bmatrix}
0.393 & 0.259 & 0.183 & 0.100 & 0.064 \\
0.289 & 0.235 & 0.215 & 0.152 & 0.109 \\
0.209 & 0.193 & 0.219 & 0.202 & 0.177 \\
0.147 & 0.162 & 0.204 & 0.232 & 0.254 \\
0.112 & 0.120 & 0.172 & 0.232 & 0.364
\end{bmatrix}
\end{array}
$$

In this matrix, $a_{ij}$ is the probability that a child ends up in the national quintile $j$ given their parent is in the national quintile $i$. Thus, $a_{15} = 0.064$ is our American dream value. We did this same procedure for all fifty states. Now that we have a matrix for each state, we can use Markov chains to predict what will happen in future generations. Again, let's look at Delaware. To find out what will happen in three generations, we must take the original Delaware matrix, call it $P$, and raise it to the third power. That is, $P_{3rd\,Generation} = P^3$, the resulting matrix is shown below.

$$
\text{Parent}
\begin{array}{c}
\text{Child} \\
\begin{bmatrix}
0.257 & 0.207 & 0.199 & 0.170 & 0.166 \\
0.248 & 0.202 & 0.199 & 0.175 & 0.176 \\
0.238 & 0.197 & 0.198 & 0.180 & 0.187 \\
0.229 & 0.193 & 0.197 & 0.197 & 0.196 \\
0.221 & 0.189 & 0.197 & 0.187 & 0.205
\end{bmatrix}
\end{array}
$$

Now our American dream value across 3 generations is, $a_{15} = 0.166$. We will now go one step further. If we take our original matrix, $P$, and raise it to the 100th power, this will show us where the probabilities will level off. So, for $P_{100th\,Generation} = P^{100}$, here is the resulting matrix.

$$\text{Parent} \begin{bmatrix} 0.232 & 0.1920 & 0.192 & 0.173 & 0.179 \\ 0.233 & 0.1920 & 0.192 & 0.173 & 0.179 \\ 0.233 & 0.1920 & 0.192 & 0.173 & 0.179 \\ 0.232 & 0.1920 & 0.192 & 0.173 & 0.179 \\ 0.233 & 0.1920 & 0.192 & 0.173 & 0.179 \end{bmatrix} \overset{\text{Child}}{}$$

Here we see our American dream value across 100 generations is $a_{15} = 0.179$. This process to generate these three matrices was done for all 50 states so that we can obtain American dream values for every state. We then used these values to create colored maps of the United States as shown in the results tab.

# 4  Results

## 4.1  The American Dream

The results of achieving the American dream for all 50 states is shown in Table 1 below. As you can see, people who live in North Dakota, Wyoming, and Alaska have the greatest chances of achieving the American dream while South Carolina, Georgia, and North Carolina have the worst chances.

| 1 | ND | 0.190 | 18 | KS | 0.0958 | 35 | KY | 0.0700 |
|---|----|-------|----|----|--------|----|----|--------|
| 2 | WY | 0.161 | 19 | NY | 0.0937 | 36 | IN | 0.0686 |
| 3 | AK | 0.131 | 20 | OK | 0.0914 | 37 | LA | 0.0685 |
| 4 | SD | 0.124 | 21 | TX | 0.0884 | 38 | MD | 0.0657 |
| 5 | IA | 0.123 | 22 | NV | 0.0881 | 39 | MO | 0.0652 |
| 6 | UT | 0.120 | 23 | ID | 0.0878 | 40 | DE | 0.0644 |
| 7 | MT | 0.117 | 24 | OR | 0.0867 | 41 | VA | 0.0633 |
| 8 | NE | 0.109 | 25 | WI | 0.0863 | 42 | FL | 0.0615 |
| 9 | NJ | 0.106 | 26 | VT | 0.0861 | 43 | MI | 0.0606 |
| 10 | MN | 0.105 | 27 | PA | 0.0857 | 44 | OH | 0.0558 |
| 11 | WV | 0.105 | 28 | NM | 0.0841 | 45 | AL | 0.0529 |
| 12 | WA | 0.104 | 29 | RI | 0.0821 | 46 | TN | 0.0525 |
| 13 | CA | 0.102 | 30 | ME | 0.0806 | 47 | MS | 0.0466 |
| 14 | MA | 0.0993 | 31 | CT | 0.0786 | 48 | NC | 0.0464 |
| 15 | NH | 0.0992 | 32 | AZ | 0.0722 | 49 | GA | 0.0429 |
| 16 | HI | 0.0976 | 33 | AR | 0.0705 | 50 | SC | 0.0408 |
| 17 | CO | 0.0973 | 34 | IL | 0.0701 | | | |

Table 2: American Dream values for all 50 states

The next table shows our predicted American Dream values in three generations. We see that the same group of states (North Midwestern States) stay at the highest values and the same group of states (Southern States) stay at the lowest values. It is also important to note that the values for all states improved significantly from the current American Dream values. This is to be expected just from the inherent Markov properties.

| 1 | ND | 0.344 | 18 | CT | 0.204 | 35 | ME | 0.171 |
|---|----|-------|----|----|-------|----|----|-------|
| 2 | WY | 0.292 | 19 | TX | 0.201 | 36 | VA | 0.168 |
| 3 | IA | 0.276 | 20 | RI | 0.201 | 37 | HI | 0.167 |
| 4 | SD | 0.260 | 21 | LA | 0.200 | 38 | OR | 0.166 |
| 5 | NE | 0.251 | 22 | NH | 0.199 | 39 | DE | 0.166 |
| 6 | MN | 0.244 | 23 | WA | 0.195 | 40 | OH | 0.163 |
| 7 | NJ | 0.232 | 24 | CO | 0.189 | 41 | AL | 0.161 |
| 8 | WV | 0.231 | 25 | VT | 0.189 | 42 | NV | 0.159 |
| 9 | KS | 0.221 | 26 | MD | 0.186 | 43 | MS | 0.153 |
| 10 | AK | 0.219 | 27 | IL | 0.186 | 44 | MI | 0.150 |
| 11 | PA | 0.219 | 28 | AR | 0.181 | 45 | AZ | 0.149 |
| 12 | MT | 0.217 | 29 | CA | 0.181 | 46 | FL | 0.144 |
| 13 | MA | 0.216 | 30 | IN | 0.177 | 47 | TN | 0.140 |
| 14 | OK | 0.215 | 31 | MO | 0.175 | 48 | NC | 0.131 |
| 15 | WI | 0.212 | 32 | ID | 0.174 | 49 | SC | 0.125 |
| 16 | UT | 0.209 | 33 | KY | 0.173 | 50 | GA | 0.123 |
| 17 | NY | 0.208 | 34 | NM | 0.173 | | | |

Table 3: American Dream values over 3 generations

The below table shows our final prediction for American Dream values, 100 generations from now. We once again observe the same trends in this table. The North Midwestern states still obtain the best scores and the Southern states still obtain the worst scores. Also, we once again see all values improve across the board, although by not nearly as much as before. At this point, the changes in our American Dream values have decelerated to a point where they are as close to stationary as we may reasonably expect. For example, if we look at Delaware, the difference between its current American Dream value and its American Dream value 3 generations down the road is $|0.166 - 0.0644| = 0.1016$. If we advance 97 generations into the future, we will see this difference drop to almost zero. As expected, the difference between Delaware's American Dream value 97 generations in the future and its American Dream value 100 generations in the future is $|0.178 - 0.179| = 0.001$.

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | ND | 0.365 | | 18 | LA | 0.217 | | 35 | ID | 0.179 |
| 2 | IA | 0.303 | | 19 | TX | 0.215 | | 36 | AL | 0.179 |
| 3 | WY | 0.289 | | 20 | NY | 0.213 | | 37 | DE | 0.178 |
| 4 | SD | 0.271 | | 21 | MD | 0.212 | | 38 | OR | 0.178 |
| 5 | NE | 0.268 | | 22 | NH | 0.209 | | 39 | NM | 0.178 |
| 6 | MN | 0.257 | | 23 | RI | 0.207 | | 40 | ME | 0.176 |
| 7 | NJ | 0.254 | | 24 | WA | 0.206 | | 41 | MS | 0.174 |
| 8 | WV | 0.251 | | 25 | VT | 0.202 | | 42 | HI | 0.173 |
| 9 | PA | 0.238 | | 26 | IL | 0.198 | | 43 | NV | 0.169 |
| 10 | CT | 0.234 | | 27 | CO | 0.198 | | 44 | MI | 0.162 |
| 11 | MT | 0.231 | | 28 | AR | 0.196 | | 45 | TN | 0.156 |
| 12 | KS | 0.231 | | 29 | KY | 0.192 | | 46 | FL | 0.153 |
| 13 | MA | 0.226 | | 30 | VA | 0.191 | | 47 | AZ | 0.150 |
| 14 | UT | 0.222 | | 31 | IN | 0.190 | | 48 | NC | 0.145 |
| 15 | OK | 0.222 | | 32 | OH | 0.190 | | 49 | SC | 0.143 |
| 16 | AK | 0.221 | | 33 | CA | 0.189 | | 50 | GA | 0.139 |
| 17 | WI | 0.220 | | 34 | MO | 0.188 | | | | |

Table 4: American Dream values over 100 generations

Below we include all three maps that correlate with the data tables above. Notice the lightness of the North Midwestern states and the darkness of the Southern states throughout all three maps. Also, notice the overall lightening of the entire map over time.
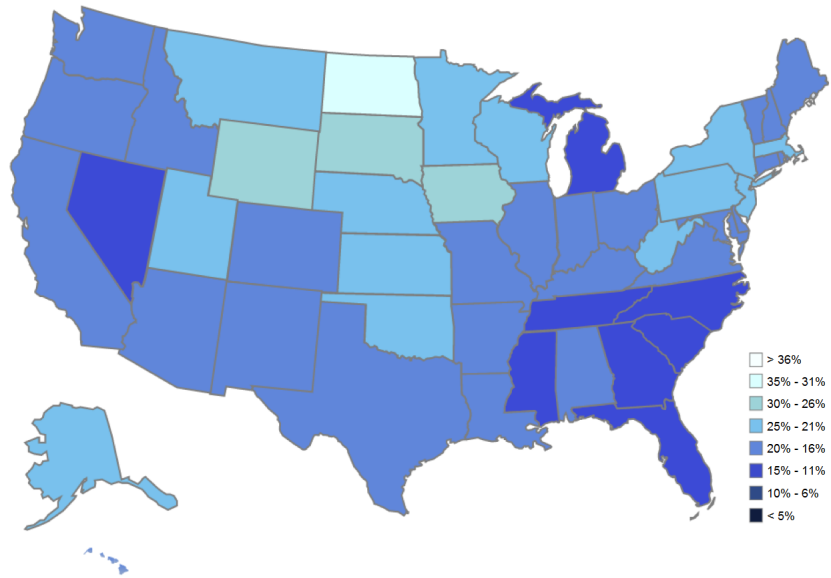


Figure 1: The current American Dream.

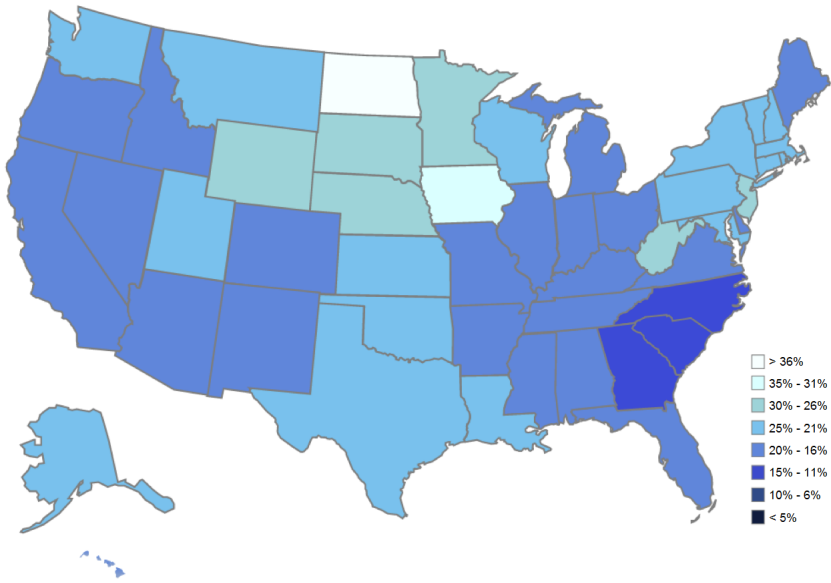Figure 2: Predicted 3rd generation American Dream.



Figure 3: Predicted 100th generation American Dream.

After analyzing this data and looking at our maps, we can see a definite trend in states with the best and worst American Dream values. Our natural

question is: What makes these states good and bad for economic mobility? In the next section we analyze various factors we hypothesize might be contributing to the difference between different areas of the United States. We look at race, education, religion, population density, and income disparity.

## 4.2   Proportion of African-Americans within each state

### 4.2.1   Map



Figure 4: Proportion of African-Americans by state

### 4.2.2   Analysis

There is a noticeable correlation between the proportion of African-Americans within each CZ and the American Dream value assigned to a state. What we see is that the proportion of African-Americans in the North Midwestern states is very low, most staying at $< 1$ percent and a couple making it into the 2 to 5 percent range. This is the lowest proportion in the whole country except for the New England area (specifically Maine, Vermont, and New Hampshire). We also see that the proportion of African-Americans is highest in the Southern States, with the average Southern state having 21 to 25 percent, and the highest proportion being greater than 31 percent in Mississippi and Louisiana.

These proportions correlate moderately with the American Dream values. The states with the highest American Dream values have the lowest proportion of African-Americans within each CZ, and the states with the lowest American Dream values have the highest proportion of African-Americans within each CZ. It is highly likely that the proportion of African-Americans in a stage correlates

with that state's American Dream score. We will analyze this relationship more in depth in our conclusion.

## 4.3   Dollars spent per K-12 Public School student
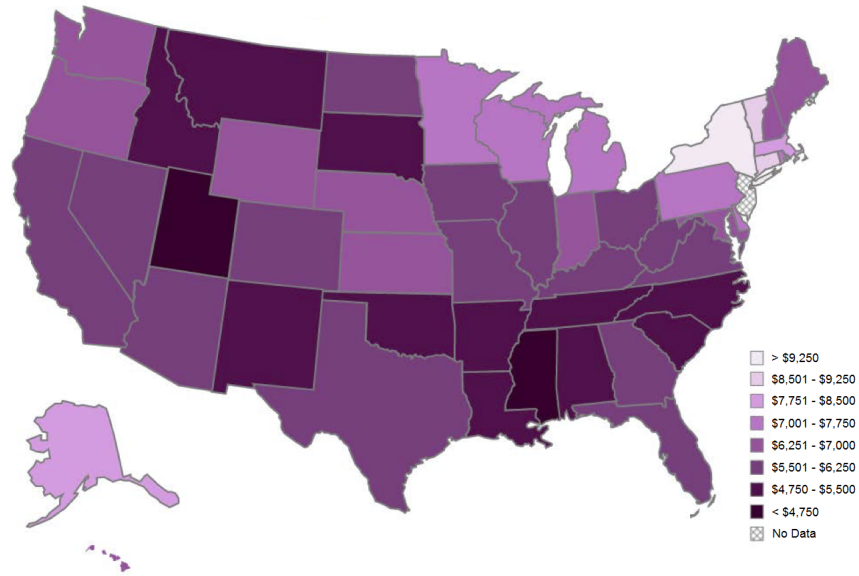
### 4.3.1   Map



Figure 5: Dollars spent per K-12 Public School student by state

### 4.3.2   Analysis

We observe that in the North Midwestern states we have a range of spending per student in K-12 public schools. We have a range from very low spending to an upper medium spending on education. In the Southern states we see slightly more consistent values, with most Southern states spending in the lower range.

It is interesting to note that both areas we are interested in include states that have the lowest spending per student in K-12 Public schools. Looking solely at this, we can find no correlation between spending on public education and the American Dream score of a state.

While from our data specifically we cannot come to any conclusions, it is important to note that there are other education factors that we have not considered here, including college attendance and percentage of children that attend private schools. These other education factors may have an impact on the American Dream values of a state, but we are unable to make a conclusion about education based on our limited data.
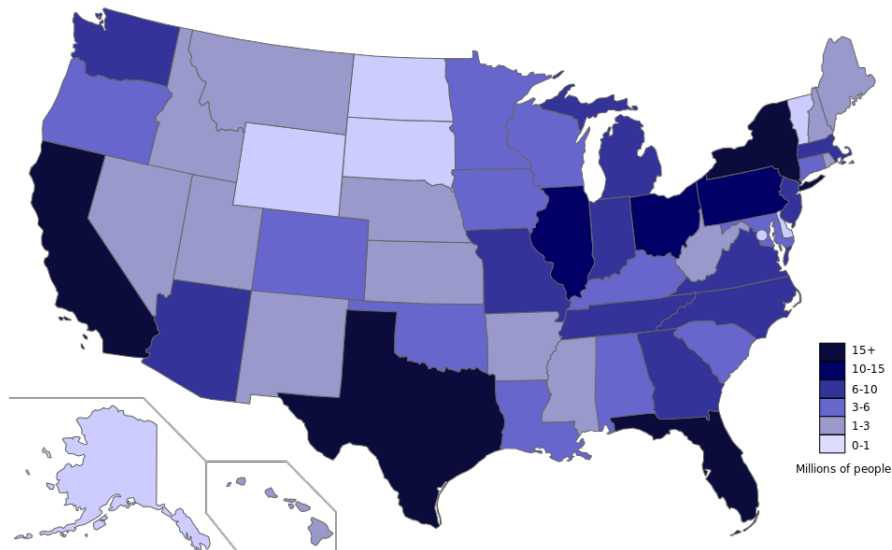
## 4.4 Proportion of people who identify as religious

### 4.4.1 Map



Figure 6: Proportion of religious people by state

### 4.4.2 Analysis

Both in the North Midwest and the South we observe that there is a wide variation in the proportion of people who identify as religious. Since there is such a wide variation, there is little to no correlation in the data and we must conclude that there is no relationship between religion and the American Dream values.

## 4.5 Population density



Figure 7: Population density by state [4]

### 4.5.1 Analysis

In the North Midwest, we observe that there is a rather low population density. North Dakota, South Dakota, and Wyoming have the lowest population densities, having 0-1 million people. In the Southern states, we see a much larger range in population density by state. The population density in the Southern states ranges from 1 to 10 million people.

While the North Midwest states stay consistent with each other, the Southern states do not. We conclude that perhaps a low population density allows for a better American Dream score to arise, but that a high population density does not determine either a bad or good American Dream score. It seems that having a low American Dream score is not determined by population density.

## 4.6 Income Disparity within States

### 4.6.1 Scatter plots

In this section we will look at income disparity within each state. We define income disparity as the difference between the high and low income in a state. We will use the 90th percentile and 10th percentile to define these low and high incomes within each CZ.

Below we have our first scatter plot. Here we are looking at the effect of Median Parent Income on American Dream values. Here we used the American Dream values for the present and not our predictions. We see a slight overall increase in American Dream values as we increase the Median parent income in

a state. However, we still have a large variation within this slight increase, with the corresponding American Dream values for a certain income ranging by up to 15 percent. It is also important to notice that our highest American Dream values from North Dakota and Wyoming are towards the middle of the median parent income. There is no obvious reason for North Dakota and Wyoming to have such high American Dream values. We must conclude that median parent income has little effect on American Dream values.
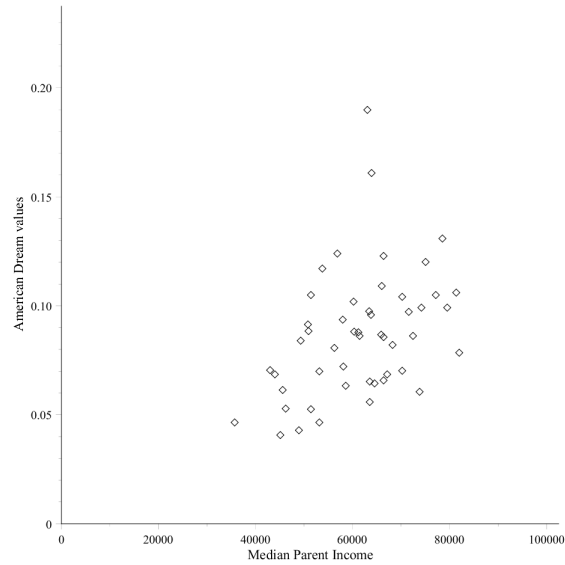


Figure 8: American Dream values vs. Median Parent Income

In our next scatter plot we look at the effect of Mean Parent Income on American Dream values. Here we have a similar pattern to that of the Median Parent Income scatter plot, except now we have even less of an increasing pattern. This scatter plot shows no correlation between income and American Dream values.
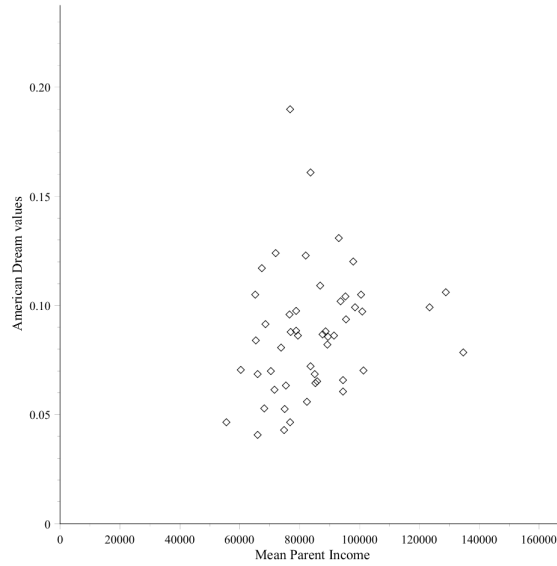
Figure 9: American Dream values vs. Mean Parent Income

Our last scatter plot compares the difference between the 90th percentile and the 10th percentile of income with the American Dream values in each state. We could not spot any trend in this scatter plot.
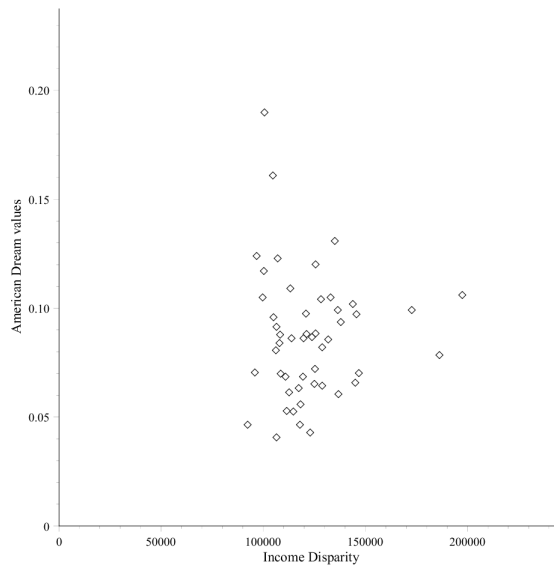


Figure 10: American Dream values vs. Income Disparity

### 4.6.2   Analysis

We had previously hypothesized that the income disparity in a state, or the difference between the highest and lowest incomes in a state, would make it easier for economic mobility to occur. However, as we see in the scatter plots, there is not a solid trend in the data. We must conclude that income disparity does not play a role in a state's American Dream score.

# 5   Conclusions and future work

While initially it would seem that education and population density should have a significant effect on upwards mobility, the only significant correlation we found was with race. Not only was race the only factor that showed a definite correlation, it showed a moderately high correlation, with our American Dream map and Race map looking almost identical. To confirm this we created a scatterplot of American Dream values vs. percentage of African Americans in each state.
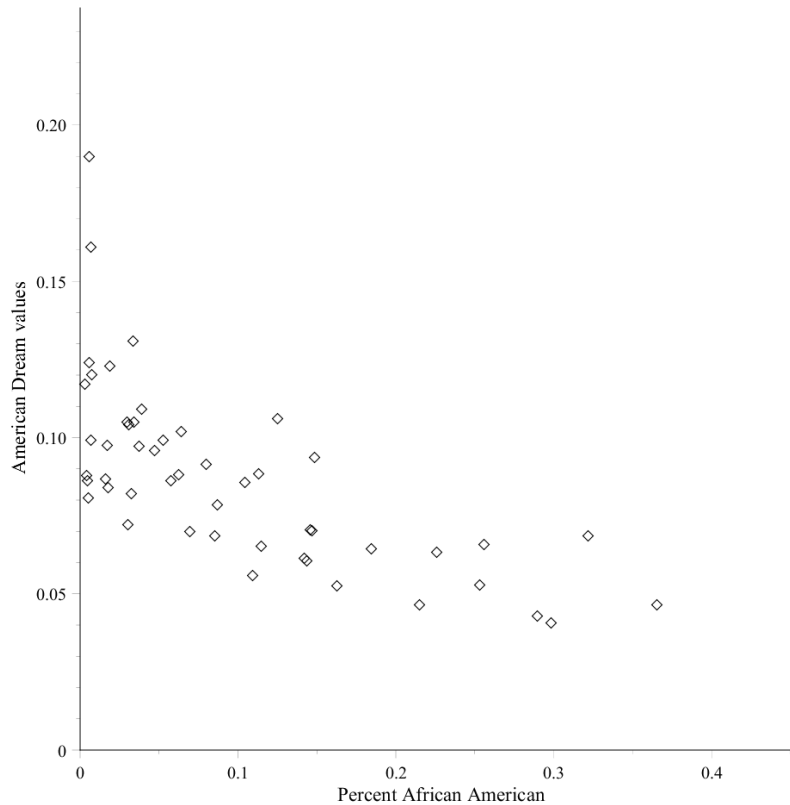


Figure 11: American Dream values vs. percentage of African Americans per state

We decided to test for the correlation coefficient for this scatter plot and we found r=-0.693. This value shows us that we have a moderately high, negative correlation.

This result fits with the results from the PEW EMP, which found that black people have a harder time moving up the ladder and are far more likely to be stuck at the bottom than other groups. In the future it is important that we further analyze why race correlates so strongly with achievement of the American Dream.

It is also important to note that there are several other factors we can check in future work, including different occupations, lifestyles, and cultures associated with different areas of the United States. There could be other factors that have a strong correlation with achievement of the American Dream that we were unable to find given our data limitations.

We would also like to note possible expansions to be made to our research in education. In our education section, we were only able to access data involving state funding of K-12 public schools. We were unable to gather data regarding other significant factors related to education, including college attendance, private school attendance, drop-out rates, and funding for college students. Any of these could shed further light on education's role in upward mobility. Another area to investigate is the intersection of education and race, and how a student may be affected by this intersection in regards to their future mobility. It was a shock to find in our studies that race has a much greater impact on economic mobility than public K-12 education, and we feel this area should be explored further.

In conclusion, there is a moderately high, negative correlation between race and upwards mobility. This conclusion is hardly a new one. We leave the task of identifying other pertinent factors to other researchers, as well as discovering reasons behind this correlation and furthermore ways in which this correlation can be lessened in the future.

# 6 Acknowledgements

# References

[1] Raj Chetty, Nathaniel Hendren, Patrick Kline, and Emmanuel Saez. Descriptive statistics by county and commuting zone: Online data table 6: Quintile-quintile transition matrices by commuting zone. 2014.

[2] Raj Chetty, Nathaniel Hendren, Patrick Kline, and Emmanuel Saez. Where is the land of opportunity? the geography of intergenerational mobility in the united states. *Journal of Economics*, 129(4):15531623, 2014.

[3] Susan K Urahn, Erin Currier, Diana Elliot, Lauren Wechsler, Denise Wilson, and Daniel Colbert. Pursuing the american dream: Economic mobility across generations, Jul 2012.

[4] Ali Zifan. The population of the states as of july 2013, based on data from united states census bureau, Nov 2014.