

Basic Info:

- **Project Name:** Demographic Analysis
- **Names:**
 - Christopher Mertin – u1010077 – cmertin@cs.utah.edu
 - SeyedMajid RasouliPichahi – u1013493 – maj.rasouli@gmail.com
 - Ashkan Bashardoust – u1011913 – u1011913@utah.edu
- **Repository:** https://github.com/cmertin/US_Stats

Overview and Motivation:

Demographic Analysis can be used for many purposes.

One example could be migration. Each year many people migrate within the US or from other countries to US for business, study, and job purposes. They try to find out which part of it is most suitable for them. Information such as percentage of highly educated people, population of the youth, and average income in that area could help a lot. Creating this kind of visualization would help these people to decide which areas suit their criteria to move.

Another example would be usage for companies. Many companies that are starting their businesses try to find the best areas for their selling market. Based on the information of average income, population of each age ranges, and the education level of people living in that area, they get a first estimate that how's their product is going to sell in that area. They can use this information for deciding whether those people are suitable for employment or not.

Other usages:

1. **Public usage:** This project could be used for finding a place to migrate, according to a specific lifestyle. For instance, a young person may have higher tendency to go to a place with younger range of people to have more fun activities and nightlife. Or to find a place having people with higher university degrees. One would prefer to live in a city, which has people with higher income rates.
2. **Journalists:** can use this project to conduct researches and use them as background for articles.
3. **Commercial companies:** can use this visualization as a basis for market research and to figure out the supply and demand ratio for their businesses.

The data we are using provides statistical analysis on certain categories around the US, will provide relevant information on the types of people in various areas of the United States. You can therefore select certain attributes or areas such that you can tailor your business to your current demographic or see to where you want to expand.

Related Works:

We have checked some population visualizations on the census.gov website, including:

1. Before and After 1940: Change in Population Density:

<https://www.census.gov/dataviz/visualizations/010/>

2. Distribution of Hispanic or Latino Population by Specific Origin:

<https://www.census.gov/dataviz/visualizations/072/>

3. Population Distribution by City Size, 1790 to 1890:

<https://www.census.gov/dataviz/visualizations/005/>

We also have checked the visualization techniques used in the final projects from the last Visualization course.

Questions:

What questions are you trying to answer?

We are trying to show how many people with those selected attribute live in a specific area.

We want to show the numerical difference of people in a certain category in each area.

How did these questions evolve over the course of the project?

What new questions did you consider in the course of your analysis?

We want to show how the population of a selection changed over the years.

Data:

The data will be provided by census.gov <http://factfinder.census.gov>

The U.S. census provides the data in easy .csv format with the info we need from 2010-2015 with the appropriate categories. The only “data cleanup” that needs to be done is to turn the columns of the data selections/attributes into percentages so that we can perform statistical analysis on those columns to give the user appropriate numbers.

For example, if a county has a population of 1 million, 450,000 of which are male, and 4,000 of the original 1 million are native American, we will change the data such that it reads 1 million for the population, 0.45 for the males, and .004 for the number of native Americans, and so forth for each of the data points. Therefore, if the user asks for the number of native American males around the US, we can say for that given county that it's approximately $1 \text{ million} * .45 * .004$ or approximately 1,800 Native American Males in that county on average. This can be done for all the attributes so the statistics can be calculated on the fly as the user selects the data.

The data that we looked at is:

- Age and Sex (population of males/females at different ages in a given region)
- Education (No High School, High School/GED, Some College, College, Graduate/Professional)
- Race
- Marital Status (Never Married, Divorced, Separated, Married, Widowed)

One problem with this data was that the census bureau did not provide the data specifically for each age. For example, the population of age and sex {18-19, 20, 21, 20-24, etc} for both males and females. To get it for each age, what we did was assume an equal distribution of people for the age ranges. For example, there were roughly 70,000 males in Alabama 18-19 years old in 2010. To make it "fair," we split the population size such that we gave 18 year old men a population size of 35,000 and the same with 19. This was a fair assumption as most of these should be a relatively equal distribution since the ranges between the years were quite small.

The rest of the data (not dealing with age and sex) was provided as a percentage of the population (at a given age), so we stored the fractional values as each column. Therefore, we can use what we stated before of taking the population size and multiplying through by each percentage value to get a rough number on the number of people of the selected attributes.

The data was parsed into two different formats, JSON and CSV. This was first output as a JSON file for each of the year for all of the states and counties, but the resulting file size was too large. Therefore, we opted to add in a CSV file as well, which greatly reduced the size of the files.

Each of the columns in the CSV file represents each combination of attribute, for example 18_M_No_HS is the percentage of the number of 18 year old men in that geographic region with No High School degree. The columns are broken up by age and gender, and there is a permutation of each of the attributes. This may need to be changed later when we try to access the data if we cannot find an efficient way to convert it to JSON for easy and quick access after reading in the file for the first time.

Exploratory Data Analysis:

We are using a US map as our major visualization to visualize our data. The user should be able to select multiple attributes such as age, gender, education level, income and race. Then the map should show color map of the number of people in that selection in each state.

Demographic Analysis

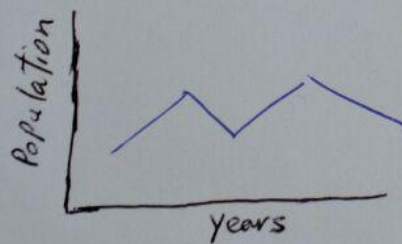
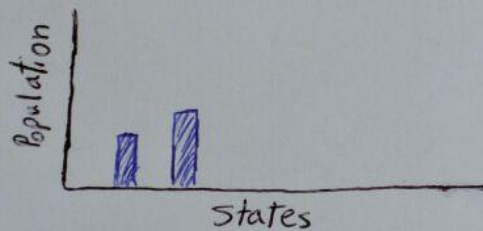
Attribute Selection

Year
...

Gender
...

Education
...

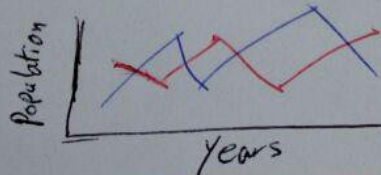
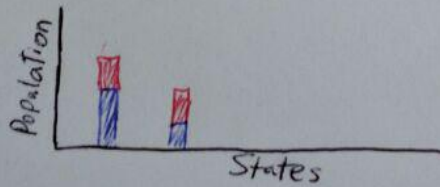
Race ...



Category Selection

category
...

Year



There will be sliding scales to select ranges, from which the data will be parsed and the results will be populated on the screen.

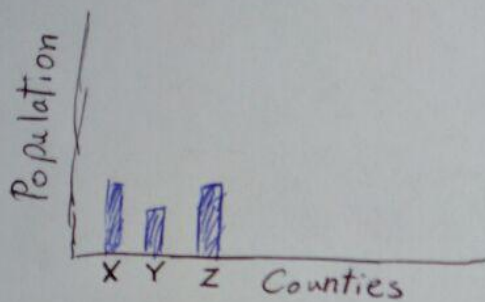
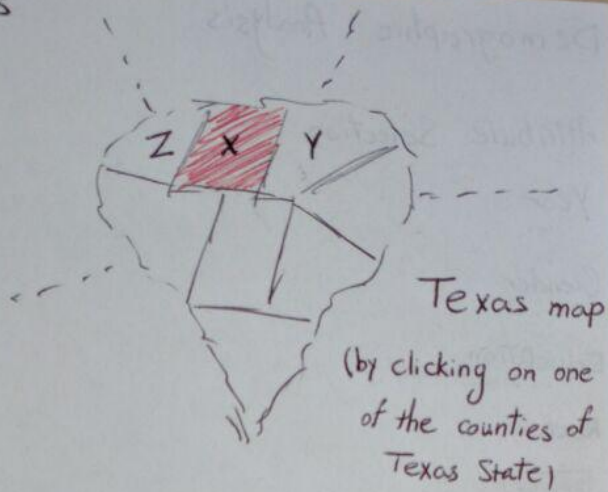
The idea is to have a color scale on the map such that it will act as a “heat map” that will show the data all around the US. Then, user can also click on a certain state and it will “zoom in” to the county level, for which it will have more precise colors as well for each individual county.

It’s necessary to mention that we have 3 choices for selecting the area:

1. The default is the whole US. Bar chart compares all the states and line chart uses the population of people in the whole country.
2. By clicking a state the map zooms into that state, showing all the counties of it. Bar chart compares all counties and line chart uses the population of people in that state.
3. By selecting a county, bar chart doesn’t change (comparing all counties) and line chart uses the population of people in that county.

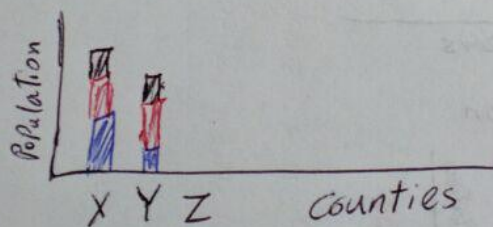
Demographic Analysis

Attribute Selection



Population change
for county X

Category Selection



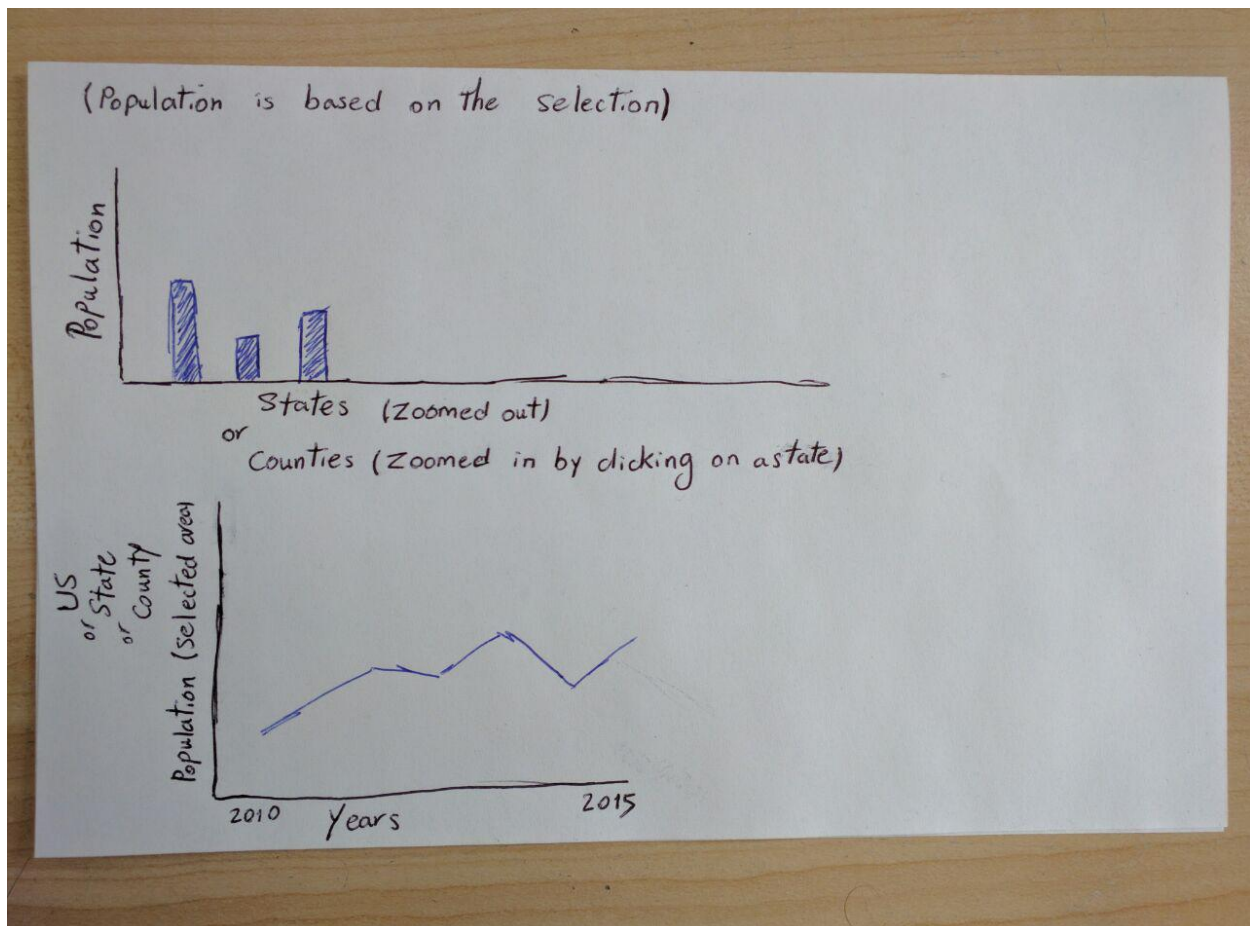
For better analyzing our data, we are using 2 different methods:

1. First method is that user chooses a selection of all categories. Then we use this information for visualization.

We visualize these charts below of our map.

There will be a line chart with one line representing the population of people with those selection attributes in the selected area over the years.

The other chart would be a bar chart with a bar for each state/county (based on the area chosen), which shows the population of people with that selected attributes in each state/county. It helps to compare this population of all states/counties together.



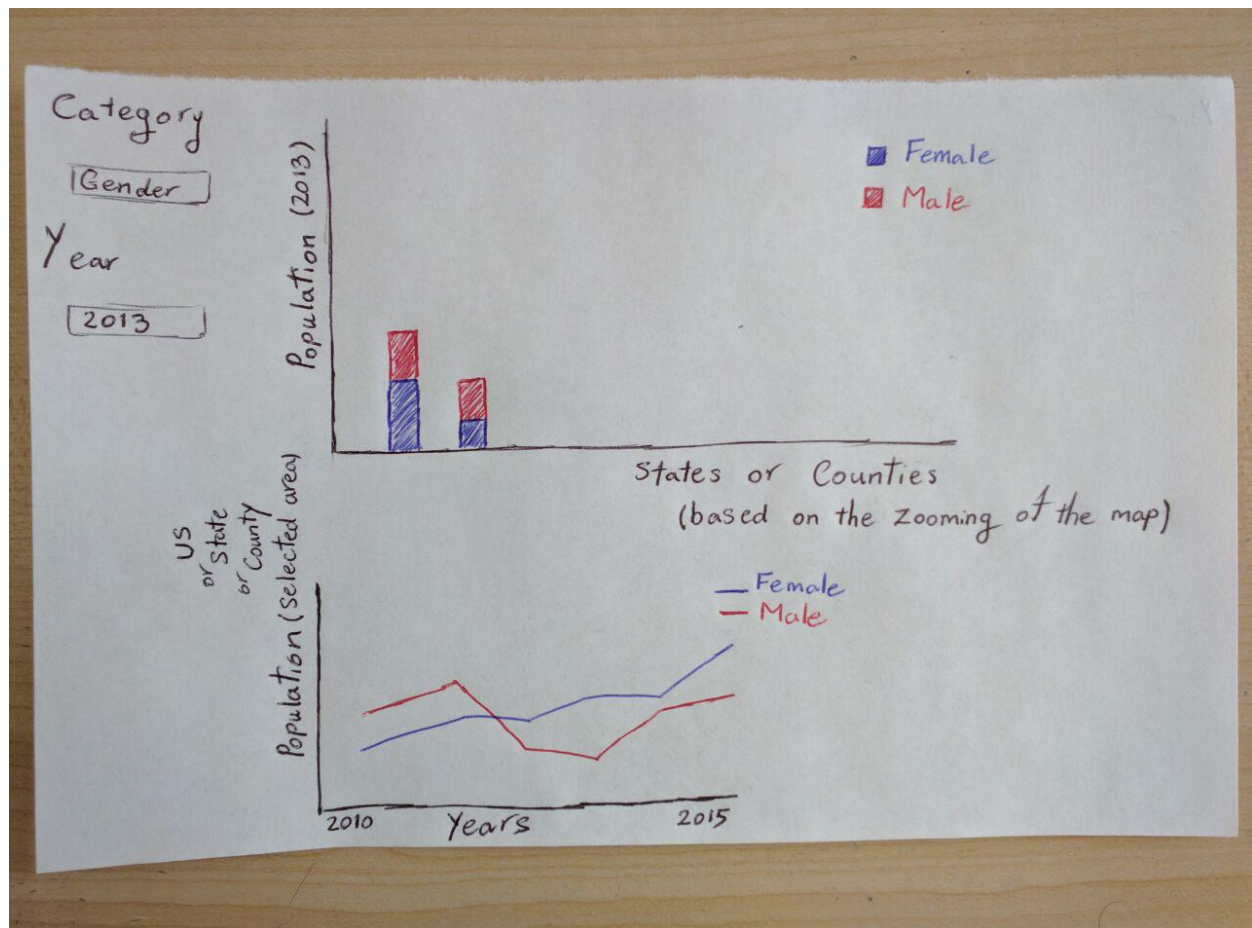
2. Second method is that user chooses a category. Then we use it for visualizing and comparing subcategories.

We visualize these charts at the bottom of the page.

There will be a line chart with a line for each subcategory, which shows the change of the population of people in that subcategory in selected area over the years.

The other chart would be a stacked bar chart with a bar for each state/county (based on the area chosen) which shows the population of all subcategory in each

state/county stacked on each other. It helps us compare the subcategories of an individual state/county, also compare each state to the other one.



Design Evolution:

We thought about using donut chart, but since there are 50 states, it would be hard to read the chart.

We were considering using as much as the subcategories of each category bar chart for each state. Considering there are 50 states, and we may have 8 bars for each states, we should visualize 400 bar charts. So we decided to change that to stacked bar chart.

We decided to have some features that are necessary for our project:

Some features that are needed for this project to be successful are a good color scale such that the data can be fairly represented and easily understood. On top of this, we need to make sure that it is customizable enough that the users can get the results that they want and that it will represent the country appropriately. While some of these attributes won't be completely

independent (for example income and race), this should provide a good enough approximation to the data.

The other features that we could use for our project:

We could potentially let the user decide the color scale they would like to use and also the type of chart they would like to see for the info on the given county (donut chart, stacked bar chart, pie chart, etc.). This is not required to complete the project, but some user customization on this aspect would make it better for the user in some aspects.

Group Critiquing:

When we had to meet with another group, they had brought up good ideas, but the most notable was that we should be able to compare two states. In our current implementation, they mentioned that it may be difficult to compare two states with a stacked bar chart (depending on the attributes) since they're not aligned. They suggested that we add an additional chart to it so that the comparison would be easier.

Implementation:

Evaluation: