

# US Census Data Analysis

Vasu Sharma

Department of Computer  
Science  
University of Colorado, Boulder  
Boulder, Colorado  
vasu.sharma@colorado.edu

Yan Zhan

Department of Economics  
University of Colorado, Boulder  
Boulder, Colorado  
yan.zhan@colorado.edu

Sutianjie Zhou

Department of Economics  
University of Colorado Boulder  
Boulder, Colorado  
sutianjie.zhou@colorado.edu

## Problem Statement

Census data provides important information about the whole nation. Census data can provide socio-economic information for citizens of the country. IPUMS[1] is an organization that preserves and harmonizes U.S. census microdata and provides easy access to this data with enhanced documentation. In this project, we aim to focus on determining socio-economic trends based on the dataset collected by IPUMS. These trends are collected based on person-level attributes such as education, age, gender, race, employment status, marital status, labor force participation etc. Using data mining techniques, we aim to identify the socio-economic trends by analysing a person's education background and their identity group to predict their earning potential, factors correlated to higher income, change in rate of return on education and income inequality over the past few decades.

## Literature Survey

Paper using IPUMS dataset for analysis of social-economic questions have been published in many peer review journals in recent decades.

Sheng et al.[2] performed data mining on census data using CART which is a decision tree based classification model. Their paper was based on

using this model to classify inhabitants in the provinces of Chengyang and Laixi.

Federal and local governments[3] use census dataset to assess the economic wellbeing of communities, identify and assist low-income and marginalized communities, and allocate funding to different social programs.

Jaison R. et al [4] used the IPUMS U.S. dataset collection to identify the rate of return on college education. They used income after graduation as a parameter to judge rate of return and discussed the problems regarding stagnant wages and the decline of wages among people without a college degree. They concluded that return has remained high in spite of rising tuition and falling earnings because the wages of those without a college degree have also been falling, keeping the college wage premium near an all-time high while reducing the opportunity cost of going to school.

In Autor, Dorn, and Hanson (2013)[5], they use the PUMA-level (Public Use Microdata Area) data from IPUMS to explain the influence of increasing import from developing countries on the local U.S. labor market. They first generate local labor and employment indexes, then they map the increasing of U.S. imports from developing countries on each local labor market, to figure out the changing of the manufacturing employment indexes. They compare the local

employment indexes before and after the increasing shock of import from developing countries (especially after China returned to the WTO in 2001), and explain that the increasing imports from developing countries reduce the U.S. local employment in manufacturing industries. In Autor, Dorn, and Hanson (2019)[7], they also use the IPUMS local data to show the relation between local employment situation of young people and their marital status. The increasing import of U.S. from developing countries reduces the local market manufacturing jobs and changes local employment structure. As the number of local manufacturing jobs decreases, the wage gap among young people increases (the reduction of manufacturing will make the local job market polarization), and the first marriage age of young people increases because more individuals want to spend more time for a potential 'better' partner.

## **Proposed Work**

In this project, based on the dataset, we will analyze the following questions

1. What is the relation among a person's educational background, their identity group and earning potential?
2. What factors are generally found to correlate with higher income?
3. What is the rate of change in income inequality in the U.S. over the past decade?
4. What is the education gap among different racial groups?
5. What is the marriage status difference among people of similar age in different racial groups with different education levels?

To answer these questions, we will first perform data cleaning and preprocessing where we will remove redundant data or apply cleaning techniques to replace duplicate or Null values with measures of central tendency. Then we will transform the dataset by performing binning on income attributes using tax brackets and normalize the data set. After data cleaning and transformation, we will work on building the classification model to identify the income bin for an individual based on their socioeconomic background. We will start with support and confidence calculation and analysis between the interested attributes (whether to receive higher education, married or not, low income or not) and individuals' demographic characteristics (gender, age, race, etc.). Then, we will apply the accuracy measurement in this class to study the relation between education and never being married before. Bayesian classification will also be applied to study the marital status and income gap when we mark different individuals as low-income, middle-income and high-income. After we compare the dependent attributes with objects in different groups based on different attribute standards, we will use an ordinary linear regression model to check the effect of gender, race, and education level on individuals' income. As for education and marital status, we will use a logit/probit model to check the effect of education, gender, racial, and other demographic attributes (in this case, we will set each age from 25 to 45 as fixed effect).

## **Dataset**

The dataset is collected from the organization IPUMS[1] which provides consistent data with documentation. We are going to use the U.S. Census Person-Level data[8] for our project. The 2019 sample will be used for the visualization as well as for building the classification model. The size of this sample is around 806 MB and it has 3239553 rows and 23 columns. Some of these

columns contain important person-level information like race, age, gender, income, education, marital status and labor force participation. These columns are chosen from a wide range of columns and are specifically focused towards an individual rather than an entire household or family. The attributes we are interested in and will study includes gender, age, race, education level, location, occupation, marital status, wage income, and total income.

## Evaluation Methods

We plan to create a machine learning model which will act as a classifier for predicting earning potential by taking into account a person's socio-economic background. This prediction model will be evaluated against the test split of the original dataset. In the preprocessing phase, earning potential will be binned based on the income attribute, so the model will predict the income bin and it will be evaluated based on confusion matrix, precision, recall, accuracy and F1 score.

## Tools

The project work is divided into certain parts and all of those parts require specific tools. The tools for some of those parts are as follows :-

1. **Data Collection and Cleaning** :- For these parts of the project, we will make use of pandas and numpy. Pandas will be used for efficiently carrying out data cleaning and transformation tasks while numpy will be used for performing computational tasks on the data
2. **Data Visualization** :- Matplotlib will be used to visualize data in the form of histograms, bar charts and scatter plots.
3. **Classification and Analysis** :- Sklearn and keras will be used for the classification aspect of the project. PCAs or Autoencoders

will be used for performing dimensionality reduction on the dataset.

## Milestones

The project work will be divided into 3 parts and the timeline for those parts are as follows :-

1. **Data Cleaning and Visualization** :- We are targeting to do significant work on this and possibly finish this by 18th July. .
2. **Classification and Analysis** :- We are targeting to do this by 30th July.
3. **Income Disparity Analysis**:- We are hoping to finish this by 1st August.

## Milestones Completed

As part of our data mining project, we have completed the following milestones :-

1. **Data Cleaning and visualization** :- In this we have tried to show income correlation with various attributes of the person-level IPUMS[1] dataset as well as the milestone related to income disparity analysis. All the visualizations, plots along with the code are also committed to github. We have used plotting libraries like seaborn and matplotlib as well as data manipulation tools like numpy, pandas to achieve these visualizations and analysis. We have also started binning the income categories for the next milestone which is classification and analysis.
2. **Income Disparity Analysis** :- The disparity analysis is compared over a decade by plotting boxplots to see the differences in income gap over that time period.

The visualizations and disparity analysis help to identify the attributes that are probable causes of classifying the income of a person.

## Milestones ToDo

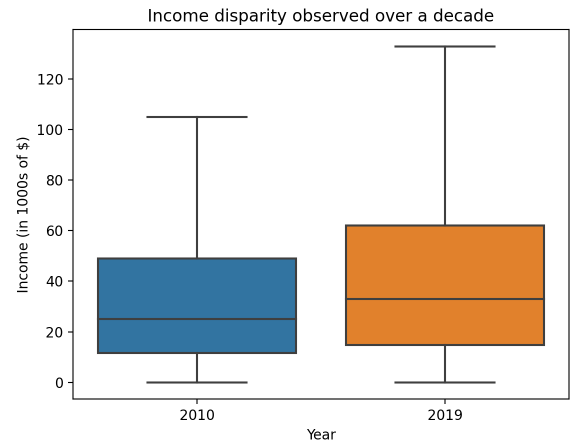
The majority of the work done around the project is around the milestones data visualization and income disparity analysis. So, the next step is to focus on the milestone :-

1. **Classification and Analysis** :- After working on the milestones related to visualization, we have tried to identify some of the key factors that are responsible for determining the income bin at a personal level. We are thinking of going ahead with a decision tree classifier as we have a key set of attributes to focus on and based on the information gain we are going to decide the order of influence of those attributes on the earning potential of a person. We also have to work on preprocessing the data for the classification task. The model will be evaluated based on the confusion matrix. If time permits, we will go ahead with another classifier i.e. a neural network classifier to compare the results that we achieved with a decision tree classifier.

## Results So Far

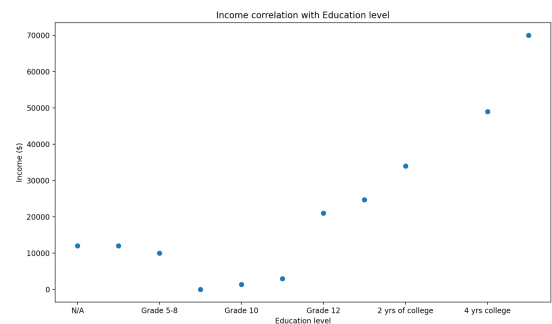
We have tried to obtain multiple visualizations of the IPUMS[1] dataset to achieve a better understanding and representation of the data. These visualizations make the job easier when it comes to milestone tasks like classification because we know about the major attributes influencing the classification variable which is the income bin of a person.

We also wanted to see how the income gap has increased over a long period of time. Income gap can explain some of the most important economic and social issues prevalent in a country. We have plotted the income gap using a boxplot which is observed over a decade.



From the boxplot of income in 2010 and 2019, we can see that the wage-gap in 2019 is larger than that in 2010. From the boxplot we can also observe that the median income of individuals hasn't increased by a significant percentage but the income among the upper middle class to rich individuals has seen a significant growth. This is one of the reasons for the ever increasing income gap in the country.

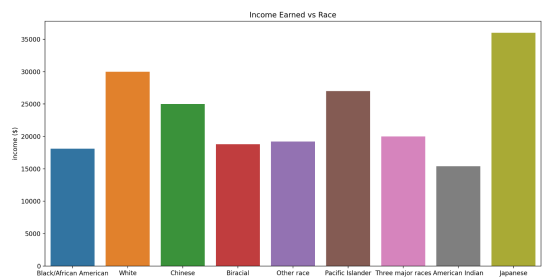
0



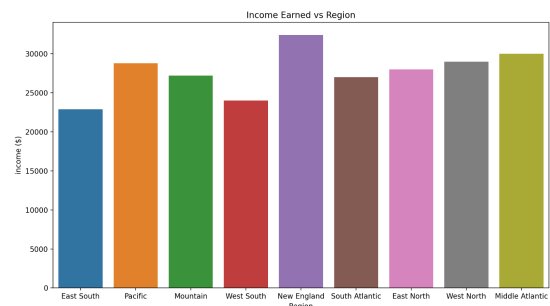
We wanted to see some correlations between income and other person level attributes in the IPUMS[1] dataset.

We first decided to do a scatter plot between income and educational background of a person. From the above visualization, we can see that income is positively correlated with the education level of a person. Another interesting fact that we observed is that having a college degree makes a significant difference in the earning potential of an individual.

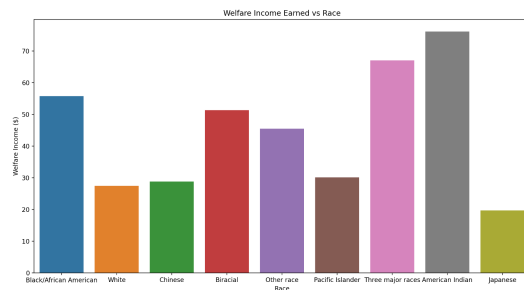
First step is to divide all individuals between 25 and 45 into 5 groups according to their education level: high-school dropout, high-school graduate, some college, college graduate, and post college. All income numbers are in dollars. The mean wage income of high-school dropout individuals is around 17000, that of high-school graduates is around 26000, while that of college graduates is around 60000 and lastly that of post graduates is around 85000. We can infer from these results that education plays a significant role in income at a personal level. This shows that a college degree is really valuable. In our data, the mean income of high school dropout female is around 11000 compared with high school dropout male which is around 21000, that of high school graduate female is around 18000 compared with high school graduate male which is around 32200, that of a graduate female is around 46000 compared with graduate male which is around 76000. We can infer from these observations that gender wage gap exists at each step of education level, especially at the higher education level where it tends to be much larger. The gender wage gap at high school dropout level is around 10000, while that at high school graduate level is around 14000 and finally at the college graduate level is around 30000. Along with the gender wage gap, we also tried to figure out a potential racial wage gap that might exist. To figure out the details about the potential racial wage gap, we generated dummy variables of white and black and observed that the mean income of black high school dropout individuals is around 10000 compared with that of white high school dropout individuals which is around 17855.56, that of black high school graduates is around 19000 compared with white high school graduates which is around 28000. These observations do hint at a potential racial wage gap prevalent in the country.



We plotted some other attributes against the total income of a person. One of them being the race of an individual. We tried to see how income is visualized among the various communities in the US. This is important because it indicates how different social and cultural communities are faring in the economic progression. These stats and visualizations help state and federal governments to create social programs that are aimed at helping certain communities which are in dire need of those assistance programs.



Location is one of the key attributes in determining the earning potential of an individual and that can be seen from the above visualization as well. The above distribution is around the median income earned by person and distributed by region. The New England region seems to have the highest median income earners among all the regions based on the IPUMS[1] dataset of year 2019.



People who are usually in the low income brackets of a country are one of the key portions of the total number of people on welfare and assistance programs. In our dataset, we utilized the variable, INCWELFR that indicates the amount of income obtained over a year from welfare or social assistance programs. We wanted to see the relationship between this particular variable and the race of an individual. This information tells about some specific communities that generally require good types of assistance programs either at the federal level or the state level. As is apparent from the above visualization, on average American Indians usually receive the most amount of income through welfare programs, followed by people who classify themselves as part of three major races and then immediately followed by Black Americans. The INCWELFR is also one of the key contributors in determining the income bin of an individual since people who generally receive more welfare from the government tend to have incomes in the lower brackets of the income split.

## REFERENCES

- [1] Steven Ruggles, Sarah Flood, Sophia Foster, Ronald Goeken, Jose Pacas, Megan Schouweiler and Matthew Sobek. IPUMS USA: Version 11.0 [dataset]. Minneapolis, MN: IPUMS, 2021. <https://doi.org/10.18128/D010.V11.0>
- [2] Bin Sheng and Sun Gengxin, "Data Mining in census data with CART," *2010 3rd International Conference on Advanced Computer Theory and Engineering(ICACTE)*,

2010, pp. V3-260-V3-264, doi: 10.1109/ICACTE.2010.5579631.

- [3] <https://www.census.gov/econ/overview/go0100.html>
- [4] Abel, Jaison R. and Deitz, Richard, Do the Benefits of College Still Outweigh the Costs? (August 1, 2014). *Current Issues in Economics and Finance*, Vol. 20, No. 3, 2014, Available at SSRN: <https://ssrn.com/abstract=2477864>
- [5] David, H., David Dorn, and Gordon H. Hanson. "The China syndrome: Local labor market effects of import competition in the United States." *American Economic Review* 103.6 (2013): 2121-68.
- [6] Steven Ruggles, Sarah Flood, Sophia Foster, Ronald Goeken, Jose Pacas, Megan Schouweiler and Matthew Sobek. IPUMS USA: Version 11.0 [dataset]. Minneapolis, MN: IPUMS, 2021. <https://doi.org/10.18128/D010.V11.0>
- [7] Dorn, David, and Gordon Hanson. "When work disappears: Manufacturing decline and the falling marriage market value of young men." *American Economic Review: Insights* 1.2 (2019): 161-78.
- [8] <https://usa.ipums.org/usa-action/variables/group>