

# **Data Mining using Classification and Clustering techniques on the Adult Census Dataset**

**Subhasree V**

**20-12-2020**

## **1. Description of the problem and a discussion of the background**

The world we live in is controlled by the economies, which is extremely dependent on an individual's income. Cambridge defines income as "money that is earned from doing work or received from investments" (INCOME | meaning in the Cambridge English Dictionary, 2020). An individual's income is very much affected by his age, occupation and unfortunately factors like gender.

The dataset used here is originally the US Census data collected in 1994. However, all of the attributes which are factors that directly or indirectly affect the income of people, is valid even today. A study conducted on 'Annual Survey of Hours and Earnings' of the year 2016 by the UK Government gives insights about the factors that can affect earnings, which is also useful for income. Age, gender, sector, skill group etc. were found to be relevant factors ( UK Government, 2020).

Understanding what affects an individual's income is important now more than ever when the world is trying hard to mitigate the effects of the Novel Coronavirus on the economy. The dataset is chosen after understanding how important income of individuals and a stable economy is for any country, especially during the time of a crisis. Another interesting article on the wall street journal suggests that the pandemic might balance wealth distribution by decreasing income equality. (Hannon, 2020) Thus the dataset that provides vital information about income and characteristics is always a relevant problem in any society.

## **2. Data description**

### **2.1. The source of the data**

The dataset originally comes from the US Census Bureau (Bureau, 2020). The United States Census Bureau is a principal agency of the U.S. Federal Statistical System, which is responsible for producing data about the people and economy of America. (United States Census Bureau, 2020) The organization releases the Census data for the public to use. Their censuses and surveys help in informed decision making and strategy building in the United States. The dataset is available in UCI Machine Learning Repository. It is the data that was extracted by Barry Becker from the 1994 Census database. The donors of this dataset are Ronny Kohavi and Barry Becker of Silicon Graphics. The data used here is downloaded from Kaggle which is named as Adult income dataset and this data is made available in Kaggle from UCI repository.

Extraction was performed by Barry Becker using some condition such as ((AAGE>16) && (AGI>100) && (AFNLWGT>1)&& (HRSWK>0)). This ensures that the dataset doesn't have lots of unclean or noisy records. (UCI Machine Learning Repository: Adult Data Set, 2020)

### **2.1.1. The agencies working with the data**

The Adult Census Dataset is cited widely and is used in various researches worldwide. The dataset was originally used by the US Census Bureau. This dataset was made publicly available and after the extraction performed by Barry Becker, the dataset continues to be used by various researchers at several academic institutions or organizations worldwide in order to study classification algorithms or to gain valuable insights. Some of the notable agencies or researchers working with the data are:

#### **“Privacy Preserving OLAP” – Paper presented on SIGMOD Conference. 2005.**

This research paper was presented on the Proceedings of the 2005 ACM SIGMOD International conference on Management of data. The authors of the paper are members of IBM Almaden and Stanford University. The paper presents the concept of data perturbation to protect from data breaches. This paper cited the Adult Census Dataset from the UCI Machine Learning Repository. The dataset was used for conducting experiments for the empirical evaluation of their algorithms.

#### **“Scaling Up the Accuracy of Naive-Bayes Classifiers: a Decision-Tree Hybrid”- Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, 1996**

The dataset was utilised in this very famous research paper that was presented in the second International Conference on Knowledge Discovery and Data Mining by Ron Kohavi. This paper introduced the Naïve Bayes Hybrid Algorithm, the NBTree which then became very popular amongst data scientists. (PDF) A Statistical Approach to Adult Census Income Level Prediction, 2020)The algorithm was performed on many many datasets available on the UCI repository, mainly the Adult Census Dataset.

### **2.2. The intended use of the data**

The data was collected and prepared to understand the factors/attributes that are associated with the income of an individual and several other characteristics that relate to and individual's socio-economic life. The dataset is mainly used for Classification and Clustering Tasks. The classification task is to determine whether or not a person makes over 50K US Dollars a year. This information can be obtained by selecting the useful features and applying classification algorithms on them. The dataset is widely used and cited, hence the maximum achieved accuracy of the classification task is also available for reference. It is also useful to identify clusters among the population which helps us better understand the clusters of people and their characteristics that affect their income, even today.

Gender discrimination in wages is still an issue faced by countries worldwide. The dataset can also be used to understand if gender affected income in the US in the 1990s. According to the data released by the U.S Department of Labor and Fortune, women make up 63 percent of workers earning the federal minimum wage. Women are in the lowest rung of the income ladder even after their significant presence in health care and food industry. (Gender Economic Inequality - Inequality.org, 2020) This clearly shows a disparity which has to be studied. Thus, the dataset has much greater significance in helping us understand the social, political and economic context that narrates the story of income in any society. New interventions with more data mining and analytics can help improve policy making and to reduce the concentration of wealth in a very small percentage of the society. Identifying patterns from this dataset is very beneficial to

develop an understanding of the most crucial factors that can increase or decrease income level of an individual.

### 2.3. Attribute types of the data

The dataset is multivariate and consists of 14 attributes

Attribute Name	Description	Type	Distinct
age	Age of the individual	Numeric	67
workclass	Work class of an individual	Nominal	8
fnlwght	The weights on the Current Population Survey (CPS) files are controlled to independent estimates of the civilian noninstitutional population of the US. These are prepared monthly for use by Population Division here at the Census Bureau	Numeric	4508
education	Education of an individual	Nominal	16
education.num	Education number of an individual	Numeric	16
marital.status	Marital status of an individual. It can be Widowed, Divorced, Separated, Never-married, Married-civ-spouse, Married-spouse-absent, Married-AF-spouse	Nominal	7
occupation	Occupation of an individual	Nominal	14
relationship	Relationship of an individual. It can be Husband, Wife, Not-in-family, Unmarried,	Nominal	6

	Own-child and other relative		
race	Race of an individual. It can be White, Black, Asian-Pac-Islander, Other and Amer-Indian-Eskimo	Nominal	5
sex	Gender of the individual. It is categorised as Female and Male	Nominal	2
capital.gain	Capital gain of an individual	Numeric	79
capital.loss	Capital loss of an individual	Numeric	56
hours.per.week	Hours per week of an individual	Numeric	78
native.country	Native country of an individual	Nominal	39
income	The class attribute. Income can be either greater than 50,000 or less than or equal to 50,000.	Nominal	2

## References

1. Archive.ics.uci.edu. 2020. *UCI Machine Learning Repository: Adult Data Set*. [online] Available at: <<http://archive.ics.uci.edu/ml/datasets/adult>> [Accessed 26 April 2020].
2. Bureau, U., 2020. *Census.Gov*. [online] Census.gov. Available at: <<https://www.census.gov/>> [Accessed 26 April 2020].
3. Dictionary.cambridge.org. 2020. *INCOME | Meaning In The Cambridge English Dictionary*. [online] Available at: <<https://dictionary.cambridge.org/dictionary/english/income>> [Accessed 24 April 2020].
4. Dl.acm.org. 2020. *Privacy Preserving OLAP | Proceedings Of The 2005 ACM SIGMOD International Conference On Management Of Data*. [online] Available at: <<https://dl.acm.org/doi/10.1145/1066157.1066187>> [Accessed 26 April 2020].
5. En.wikipedia.org. 2020. *United States Census Bureau*. [online] Available at: <[https://en.wikipedia.org/wiki/United\\_States\\_Census\\_Bureau](https://en.wikipedia.org/wiki/United_States_Census_Bureau)> [Accessed 23 April 2020].
6. Hannon, P., 2020. *How The Coronavirus Might Reduce Income Inequality*. [online] WSJ. Available at: <<https://www.wsj.com/articles/how-the-coronavirus-might-reduce-income-inequality-11587304801>> [Accessed 21 April 2020].
7. Inequality.org. 2020. *Gender Economic Inequality - Inequality.Org*. [online] Available at: <<https://inequality.org/gender-inequality/>> [Accessed 25 April 2020].
8. ResearchGate. 2020. *(PDF) A Statistical Approach To Adult Census Income Level Prediction*. [online] Available at: <[https://www.researchgate.net/publication/328494313\\_A\\_Statistical\\_Approach\\_to\\_Adult\\_Census\\_Income\\_Level\\_Prediction](https://www.researchgate.net/publication/328494313_A_Statistical_Approach_to_Adult_Census_Income_Level_Prediction)> [Accessed 24 April 2020].
9. UK Government, 2020. [online] Available at: <<https://www.youtube.com/watch?v=mo2dqHbLpQo>> [Accessed 28 May 2020].
10. UK Government. 2020. [online] Available at: <<https://www.ons.gov.uk/employmentandlabourmarket/peopleinwork/earningsandworkinghours/articles/analysisoffactorsaffectingearningsusingannualsurveyofhoursandearnings/01>> [Accessed 20 April 2020].