# CLASSIFYING ADULT CENSUS DATASET USING PYTHON AND WEKA – IBM WATSON

SUBHASREE VADUKOOT

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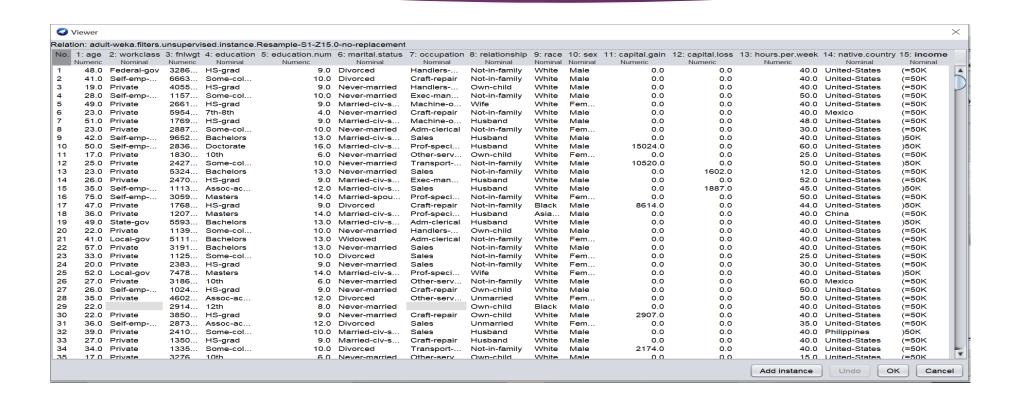
1. **Description of the problem and a discussion of the background**

▶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).
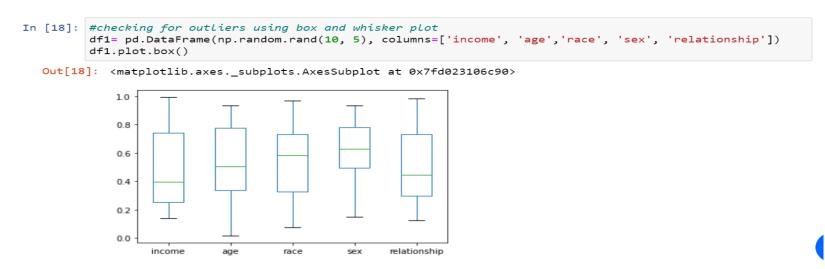
# DATA

▶ 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.

# ATTRIBUTES

# EXPLORATORY DATA ANALYSIS

BOX AND WHISKER PLOT FOR CHECKING FOR OUTLIERS

```
In [18]: #checking for outliers using box and whisker plot
         df1= pd.DataFrame(np.random.rand(10, 5), columns=['income', 'age','race', 'sex', 'relationship'])
         df1.plot.box()

Out[18]: <matplotlib.axes._subplots.AxesSubplot at 0x7fd023106c90>
```
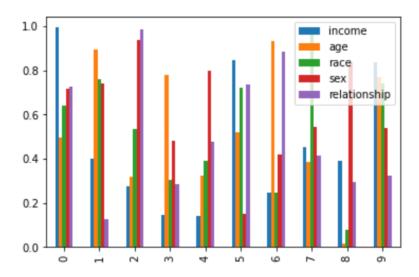


There are no outliers in the dataset

# Bar plot

# CLASSIFICATION – MACHINE LEARNING

- DECISION TREE
- KNN
- LOGISTIC REGRESSION
- SVM

# CLASSIFICATION CODE

```
In [28]: from sklearn.linear_model import LogisticRegression
         logreg = LogisticRegression()
         logreg.fit(X_train, y_train)
         print('Accuracy of Logistic regression classifier on training set: {:.2f}'
               .format(logreg.score(X_train, y_train)))
         print('Accuracy of Logistic regression classifier on test set: {:.2f}'
               .format(logreg.score(X_test, y_test)))
```

```
Accuracy of Logistic regression classifier on training set: 0.81
Accuracy of Logistic regression classifier on test set: 0.81
```

## DECISION TREES

```
In [29]: from sklearn.tree import DecisionTreeClassifier
         clf = DecisionTreeClassifier().fit(X_train, y_train)
         print('Accuracy of Decision Tree classifier on training set: {:.2f}'
               .format(clf.score(X_train, y_train)))
         print('Accuracy of Decision Tree classifier on test set: {:.2f}'
               .format(clf.score(X_test, y_test)))
```

```
Accuracy of Decision Tree classifier on training set: 0.85
```

# CONFUSION MATRIX

```
In [33]: from sklearn.metrics import classification_report
         from sklearn.metrics import confusion_matrix
         pred = clf.predict(X_test)
         print(confusion_matrix(y_test, pred))
         print(classification_report(y_test, pred))
```

```
[[5944  249]
 [1161  787]]
              precision    recall  f1-score   support

       <=50K       0.84      0.96      0.89      6193
        >50K       0.76      0.40      0.53      1948

    accuracy                           0.83      8141
   macro avg       0.80      0.68      0.71      8141
weighted avg       0.82      0.83      0.81      8141
```

# RESULTS AND CONCLUSION

▶ The best accuracy is offered by decision trees with 85% percentage.

▶ Precision for income with <=50k is 84% and greater than 50k is 76%

▶ Income is still heavily dependant on factors such as sex, marital status and other attributes.

▶ We live in an era when people should be able to earn equally when compared to their counterparts of different gender, race, region etc. This has to change.

▶ Analyses with data and machine learning can bring much needed change in this field

# 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> [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> [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].