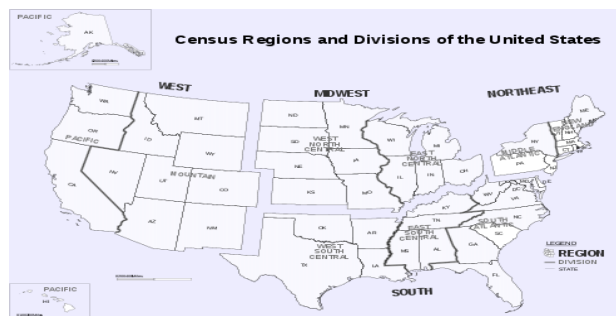Full Name          : Peter Sunny Shanthveer Markappa
Student No.        : R00208303
Subject            : Practical Machine Learning
Assignment         : 02
Data Set Used      : Adult Data Set


*********************************************************************************

# Introduction

United states census data is selected for study, data extraction was done by Barry Becker from the 1994 Census database. The United States Census Bureau (USCB) is the main agency which comes under the U.S federal statistical system (USFSS) is responsible for data collection of American peoples, which is directly under the control of the US president.  The population census is considered as a part of the constitution which has to be conducted at least once in ten years. The United State census bureau defined 9 statistical divisions in the country which are distributed in 4 regions as follows.



(Source : https://en.wikipedia.org/wiki/United_States_Census_Bureau)


**Why this study only:**

 The importance of the study,

1. To know the new area where housing and public facilities are needed.
2. To know the demographic distribution of communities in the country.
3. For better planning of cities and transportation systems.
4. Providing domestic facilities such as police stations, electricity, schools ,colleges and hospitals.
5.  For creating localised regions for conducting the election for people's choice to choose the leader in democratic country.

# Detailed Descriptions of the data.

The datasets contain 48842 records with six features continuous and eight features are nominal in nature.

**A. Continuous features are,**

> 1.age: continuous.
>
> 2. fnlwgt: continuous.
>
> 3. education-num: continuous.
>
> 4. capital-gain: continuous.
>
> 5.capital-loss: continuous.
>
> 6. hours-per-week: continuous.

**B. Nominal Features are,**

> 1. workclass: Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked.
> 2. education: Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 1st-4th, 10th, Doctorate, 5th-6th, Preschool.
> 3. relationship: Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried.
> 4. race: White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black.
> 5. marital-status: Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse.
> 6. occupation: Tech-support, Craft-repair, Other-service, Sales, Exec-managerial, Prof-specialty, Handlers-cleaners, Machine-op-inspct, Adm-clerical, Farming-fishing, Transport-moving, Priv-house-serv, Protective-serv, Armed-Forces.
> 7. sex: Female, Male.
> 8. Native-country.

**C.Target Feature**

1. income: >50K, <=50K.

**Objective:**

Building a binary classification model, i.e which can predict/classify an individual based on income that is >50k or <= 50k for the given input data.

# Report - Research

From the detailed statistical descriptions of the data revealed major pre-processing which are need to address before predictive model development and the findings are

    a. Handling Of special characters.

    b. Appropriate imputation of null value.

    c. Handling of outlier values.

    d. Handling imbalanced datasets.

    e. Optimization of Hyper-parameters.


        *The current study utilizes 'sklearn'/scikit-learn python library for processing the data.*

## A. Handling of Special Characters.

Some of the nominal features of the datasets contain special characters such as "?" in the data which will be replaced with null value. This operation can be performed using the pandas library using replace function with regular expression argument to the dataframe.


## B. Appropriate imputation of null value.

From the above procedure we will end up creating null/empty cells, these cells/values should be treated. Either we can drop the missing value records from the further analysis, because of which we end up losing the data/information and another way is imputing missing values with appropriate values such mean or mode or some constant values. In the current study missing cells/values are identified in categorical features are education,occupation and native-country. In these features imputation are carried out using SimpleImputer, a function in sklearn library with mode value is imputed.


## C. Handling of outlier values.

In the current study outliers are identified using boxplot visualization as shown in the following figure. Outliers identified in the features are educational-num, age, fnlwgt, capital-gain and capital-loss.

**Educational-num** : As shown in the first row of figure 1 box plot the observations less than the 5 are considered as an outlier; these values are replaced with the mean value.

**age**: From the below chart's second row box plot it is clear that age greater than the 80 years is considered as an outlier, these values are treated with the mean value.

**Hours-per-week** : In this variable there are more outliers identified in lower and upper bound with skewness value of 0.2. The distribution of the variable is found to be normal hence no outlier treatment is carried out.

**Fnlwgt** : From the below chart the fourth row box plot shows a large number of outliers in the data with skewness of 1.4, When keenly looked at the data all observations are found to be unique hence no treatment is carried out.

**Capital - loss and Capital gain** :In these two features behaviour is found to be more or less same with a large number of outliers in the data from plots, but approximately 97% of observations are zero only 3% observations are non-zero hence no outlier treatment was carried out.

D. **Optimization of Hyper-parameters :** GridSearchCV function is used to tune the hyper parameters for randomforestclassifier.
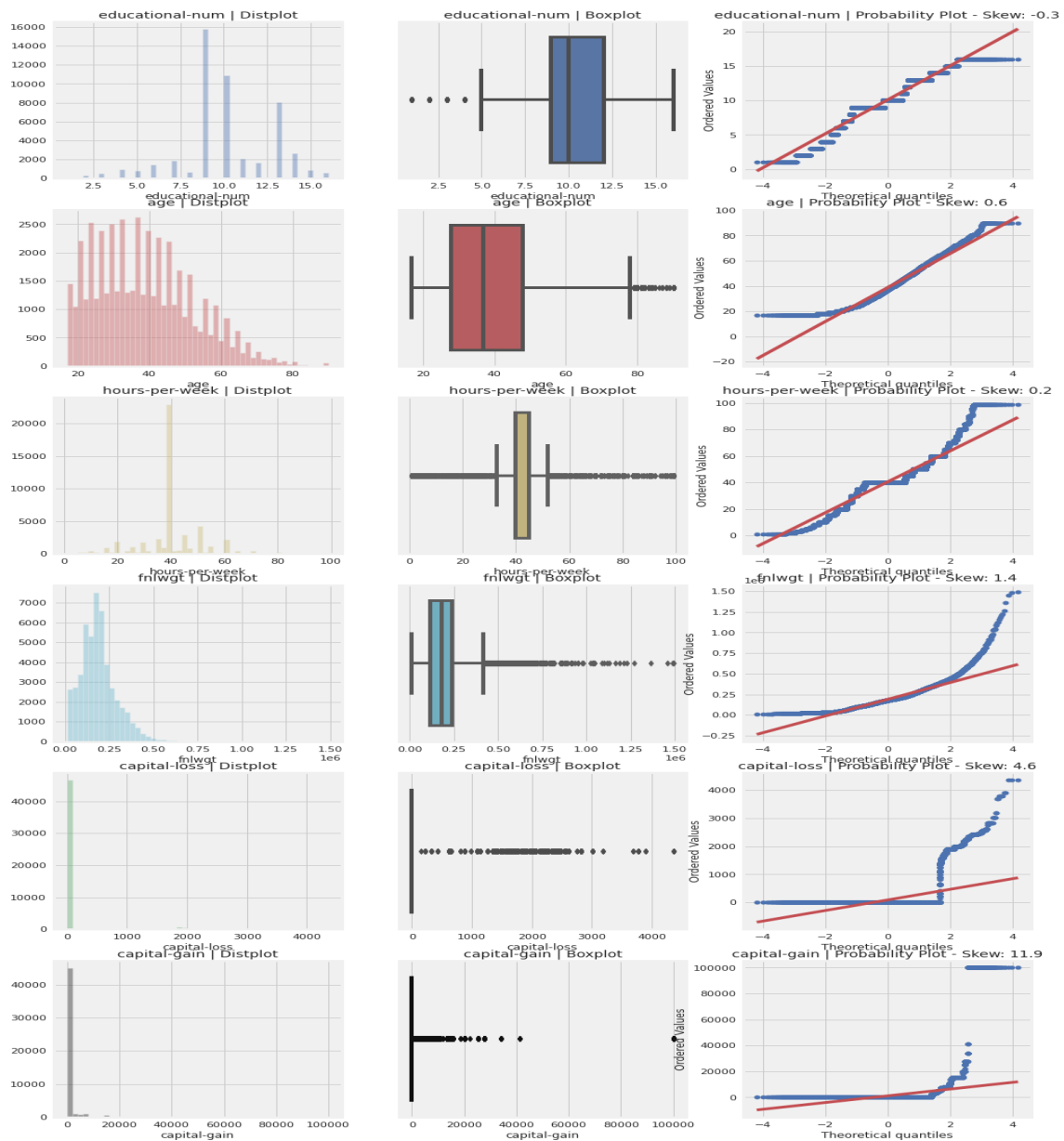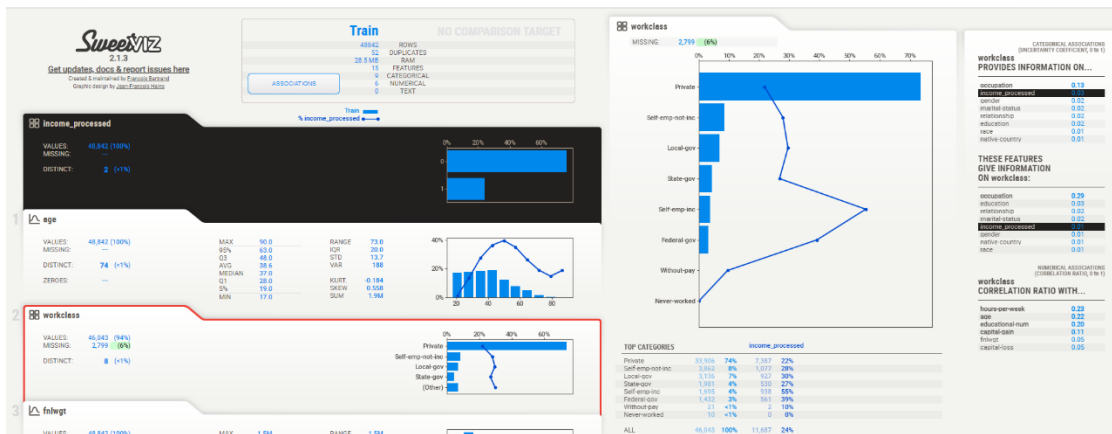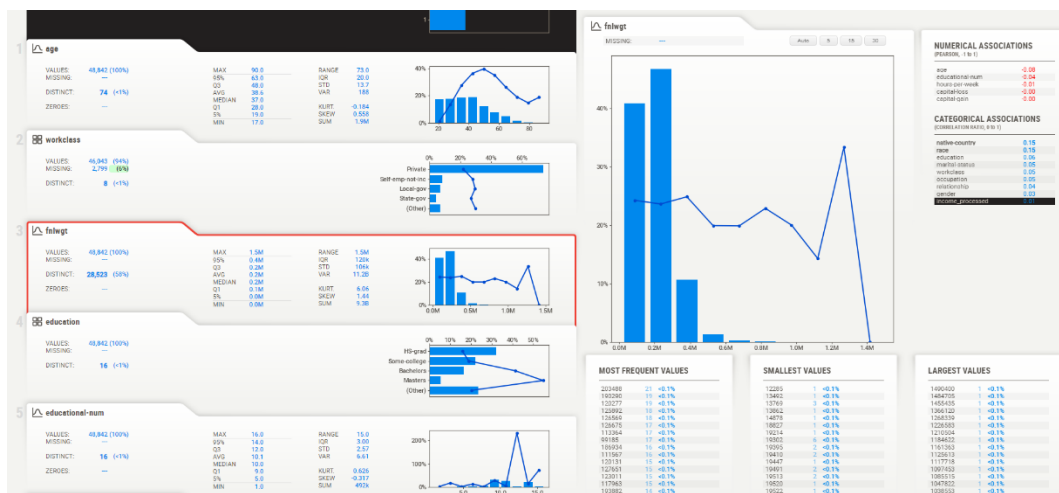


Fig:1

# Chart Created using Sweetviz
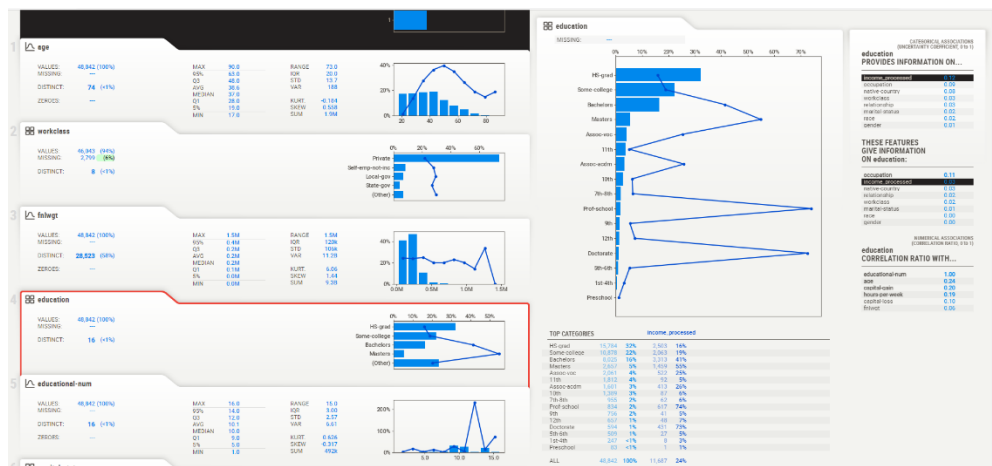


**Income_Processed**



**Age Details**



**WorkClass**

**Fnlwgt**

**Education**

**Educational-num**

## Marital-status

**marital-status**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 7 (<1%)

**occupation**
VALUES: 46,033 (94%)
MISSING: 2,809 (6%)
DISTINCT: 14 (<1%)

Prof-specialty
Craft-repair
Exec-managerial
Adm-clerical
(Other)

**relationship**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 6 (<1%)

Husband
Not-in-family
Own-child
Unmarried
(Other)

**race**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 5 (<1%)

White
Black
Asian-Pac-Islander
Amer-Indian-Eskimo
Other

**marital-status**
MISSING: —

Married-civ-spouse
Never-married
Divorced
Separated
Widowed
Married-spouse-absent
Married-AF-spouse

CATEGORICAL ASSOCIATIONS
(UNCERTAINTY COEFFICIENT, 0 to 1)
**marital-status PROVIDES INFORMATION ON...**

| | |
|---|---|
| relationship | 0.49 |
| income_processed | 0.21 |
| gender | 0.18 |
| race | 0.02 |
| occupation | 0.02 |
| workclass | 0.02 |
| native-country | 0.02 |
| education | 0.01 |

**THESE FEATURES GIVE INFORMATION ON marital-status:**

| | |
|---|---|
| relationship | 0.57 |
| gender | 0.09 |
| income_processed | 0.07 |
| occupation | 0.04 |
| education | 0.02 |
| workclass | 0.02 |
| race | 0.01 |
| native-country | 0.01 |

NUMERICAL ASSOCIATIONS
(CORRELATION RATIO, 0 to 1)
**marital-status CORRELATION RATIO WITH...**

| | |
|---|---|
| age | 0.58 |
| hours-per-week | 0.25 |
| educational-num | 0.11 |
| capital-gain | 0.08 |
| capital-loss | 0.08 |
| fnlwgt | 0.05 |

**TOP CATEGORIES** income_processed

| | | | | |
|---|---|---|---|---|
| Married-civ-spouse | 22,379 | 46% | 9,984 | 45% |
| Never-married | 16,117 | 33% | 733 | 5% |
| Divorced | 6,633 | 14% | 671 | 10% |
| Separated | 1,530 | 3% | 99 | 6% |
| Widowed | 1,518 | 3% | 128 | 8% |
| Married-spouse-absent | 628 | 1% | 58 | 9% |
| Married-AF-spouse | 37 | <1% | 14 | 38% |
| ALL | 48,842 | 100% | 11,687 | 24% |

### Marital-status

## Occupation

**occupation**
VALUES: 46,033 (94%)
MISSING: 2,809 (6%)
DISTINCT: 14 (<1%)

Prof-specialty
Craft-repair
Exec-managerial
Adm-clerical
(Other)

**relationship**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 6 (<1%)

Husband
Not-in-family
Own-child
Unmarried
(Other)

**race**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 5 (<1%)

White
Black
Asian-Pac-Islander
Amer-Indian-Eskimo
Other

**gender**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 2 (<1%)

Male
Female

**capital-gain**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 123 (<1%)
ZEROES: 44,807 (92%)

| MAX | 100k | RANGE | 100k |
|---|---|---|---|
| 95% | 5k | IQR | 0.00 |
| Q3 | 0k | STD | 7,452 |
| AVG | 1k | VAR | 55.5M |
| MEDIAN | 0k | | |
| Q1 | 0k | KURT. | 153 |
| 5% | 0k | SKEW | 11.9 |
| MIN | 0k | SUM | 52.7M |

**occupation**
MISSING: 2,809 (6%)

Prof-specialty
Craft-repair
Exec-managerial
Adm-clerical
Sales
Other-service
Machine-op-inspct
Transport-moving
Handlers-cleaners
Farming-fishing
Tech-support
Protective-serv
Priv-house-serv
Armed-Forces

CATEGORICAL ASSOCIATIONS
(UNCERTAINTY COEFFICIENT, 0 to 1)
**occupation PROVIDES INFORMATION ON...**

| | |
|---|---|
| workclass | 0.29 |
| gender | 0.16 |
| income_processed | 0.11 |
| education | 0.11 |
| relationship | 0.06 |
| marital-status | 0.04 |
| native-country | 0.04 |
| race | 0.02 |

**THESE FEATURES GIVE INFORMATION ON occupation:**

| | |
|---|---|
| workclass | 0.13 |
| education | 0.09 |
| gender | 0.04 |
| relationship | 0.03 |
| income_processed | 0.02 |
| marital-status | 0.02 |
| native-country | 0.01 |
| race | 0.01 |

NUMERICAL ASSOCIATIONS
(CORRELATION RATIO, 0 to 1)
**occupation CORRELATION RATIO WITH...**

| | |
|---|---|
| educational-num | 0.56 |
| hours-per-week | 0.31 |
| age | 0.18 |
| capital-gain | 0.12 |
| capital-loss | 0.08 |
| fnlwgt | 0.05 |

**TOP CATEGORIES** income_processed

| | | | | |
|---|---|---|---|---|
| Prof-specialty | 6,172 | 13% | 2,784 | 45% |
| Craft-repair | 6,112 | 13% | 1,383 | 23% |
| Exec-managerial | 6,086 | 13% | 2,908 | 48% |
| Adm-clerical | 5,611 | 12% | 768 | 14% |
| Sales | 5,504 | 12% | 1,475 | 27% |
| Other-service | 4,923 | 11% | 204 | 4% |
| Machine-op-inspct | 3,022 | 7% | 372 | 12% |
| Transport-moving | 2,355 | 5% | 481 | 20% |
| Handlers-cleaners | 2,072 | 5% | 138 | 7% |
| Farming-fishing | 1,490 | 3% | 173 | 12% |
| Tech-support | 1,446 | 3% | 420 | 29% |
| Protective-serv | 983 | 2% | 308 | 31% |
| Priv-house-serv | 242 | <1% | 3 | 1% |
| Armed-Forces | 15 | <1% | 5 | 33% |
| ALL | 46,033 | 100% | 11,687 | 24% |

### Occupation

## Relationship

Adm-clerical
(Other)

**relationship**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 6 (<1%)

Husband
Not-in-family
Own-child
Unmarried
(Other)

**race**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 5 (<1%)

White
Black
Asian-Pac-Islander
Amer-Indian-Eskimo
Other

**gender**
VALUES: 48,842 (100%)
MISSING: —
DISTINCT: 2 (<1%)

Male
Female

**capital-gain**
VALUES: 48,842 (100%)
MISSING: —

| MAX | 100k | RANGE | 100k |
|---|---|---|---|
| 95% | 5k | IQR | 0.00 |

**relationship**
MISSING: —

Husband
Not-in-family
Own-child
Unmarried
Wife
Other-relative

CATEGORICAL ASSOCIATIONS
(UNCERTAINTY COEFFICIENT, 0 to 1)
**relationship PROVIDES INFORMATION ON...**

| | |
|---|---|
| marital-status | 0.57 |
| gender | 0.43 |
| income_processed | 0.21 |
| occupation | 0.03 |
| race | 0.02 |
| workclass | 0.02 |
| education | 0.02 |
| native-country | 0.02 |

**THESE FEATURES GIVE INFORMATION ON relationship:**

| | |
|---|---|
| marital-status | 0.49 |
| gender | 0.18 |
| income_processed | 0.08 |
| occupation | 0.06 |
| education | 0.03 |
| workclass | 0.02 |
| race | 0.01 |
| native-country | 0.01 |

NUMERICAL ASSOCIATIONS
(CORRELATION RATIO, 0 to 1)
**relationship CORRELATION RATIO WITH...**

| | |
|---|---|
| age | 0.47 |
| hours-per-week | 0.31 |
| educational-num | 0.16 |
| capital-gain | 0.09 |
| capital-loss | 0.08 |
| fnlwgt | 0.04 |

**TOP CATEGORIES** income_processed

| | | | | |
|---|---|---|---|---|
| Husband | 19,716 | 40% | 8,846 | 45% |
| Not-in-family | 12,583 | 26% | 1,276 | 10% |
| Own-child | 7,581 | 16% | 111 | 1% |
| Unmarried | 5,125 | 10% | 309 | 6% |
| Wife | 2,331 | 5% | 1,093 | 47% |
| Other-relative | 1,506 | 3% | 52 | 3% |
| ALL | 48,842 | 100% | 11,687 | 24% |

### Relationship

**Race**

**gender**

**Capital-gain**

**Capital-loss**


**Hours-per-week**


**Native-country**

Heatmap of all Continuous Variables including target =

**Heatmap**



Violin Plot of all Continuous Variables

**Violin Plot**

**Scatter-plot**

# Methodology

In this section detailed analytical methodology will be described such as basic statistical analysis and data pre-processing steps are explained in the following,

1. Establishing a baseline.
2. Basic experimentation.

**1. Establishing a baseline: Basic descriptive** statistical analysis is carried out to know the distribution of variables using histogram, boxplot for outlier identification and probability chart to know the skewness for the continuous variables.

**2. Basic experimentation:** Data pre-processing includes data cleaning such as outlier handling, features selection, data transformations and data balancing.
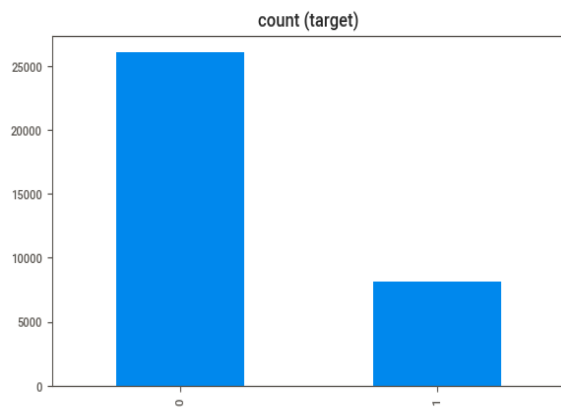
> **Data Cleaning** : Removal of symbols or special characters from the data because mathematical operations can not be carried out in such data.
>
> **Missing value:** Missing value is imputed with the most frequent value of the variable in the datasets, because missing values are identified in categorical variables imputed with the mode value.
>
> **Feature selection :** some of features are dropped from the analysis such as educational-num and fnlwgt because of low correlation with the target variables.
>
> **Data transformations :** One hot encoding transformation is carried out for the categorical variables.
>
> **Handling imbalanced datasets:** In the current study target variables(income)  with class "<=50K" are more number of observations compared to class ">50K". Data balance is carried out using "SMOTE".(Fig-2: Data Imbalance  and Fig -3: Data Balance with SMOTE).



**(Fig-2)**



**(Fig-3)**

# Evaluation and Conclusions

In the study of census data of the US found that the mean age of population was 38.64 years with the standard deviation(SD) of 13.71 years, where as minimum and maximum age of respondents are 17 and 90 years respectively.The average working hours per week is 40 hrs with a SD of 12.39 hrs. 32% of respondents from the datasets have HS-grade followed by some-college is 22%. 46% of respondents are Married-civ-spouses who have marital status, followed by Never-married 33%. 86% of respondents are white race and 10% are black race. 67% of respondents are male and 33% respondents are female. 74% of respondents are belong private workclass and 8% of respondents are self-employed.

The classification models build for the study are **LogisticRegression, LinearSVC, RandomForestClassifier, GradientBoostingClassifier, BaggingClassifier, SVC, DecisionTreeClassifier and XGBClassifier**

And results of 10 fold cross validation are presented in table-1. The accuracy values for models in test datasets are found to be

| | |
|---|---|
| LogisticRegression | : 62% |
| LinearSVC | : 75% |
| RandomForestClassifier | : 98% |
| GradientBoostingClassifier | : 88% |
| BaggingClassifier | : 98% |
| SVC | : 63% |
| DecisionTreeClassifier | : 98% |
| XGBClassifier | : 92% |

Respectively and results are saved in csv file.

```
Results for LogisticRegression:
              precision    recall  f1-score   support

           0       0.57      0.93      0.71     26043
           1       0.81      0.31      0.45     26043

    accuracy                           0.62     52086
   macro avg       0.69      0.62      0.58     52086
weighted avg       0.69      0.62      0.58     52086


Results for LinearSVC:
              precision    recall  f1-score   support

           0       0.77      0.71      0.74     26043
           1       0.73      0.79      0.76     26043

    accuracy                           0.75     52086
   macro avg       0.75      0.75      0.75     52086
weighted avg       0.75      0.75      0.75     52086


Results for RandomForestClassifier:
              precision    recall  f1-score   support

           0       0.98      0.97      0.98     26043
           1       0.98      0.98      0.98     26043

    accuracy                           0.98     52086
   macro avg       0.98      0.98      0.98     52086
weighted avg       0.98      0.98      0.98     52086
```

```
Results for GradientBoostingClassifier:
              precision    recall  f1-score   support

           0       0.90      0.86      0.88     26043
           1       0.87      0.91      0.89     26043

    accuracy                           0.88     52086
   macro avg       0.89      0.88      0.88     52086
weighted avg       0.89      0.88      0.88     52086


Results for BaggingClassifier:
              precision    recall  f1-score   support

           0       0.98      0.97      0.98     26043
           1       0.98      0.98      0.98     26043

    accuracy                           0.98     52086
   macro avg       0.98      0.98      0.98     52086
weighted avg       0.98      0.98      0.98     52086


Results for SVC:
              precision    recall  f1-score   support

           0       0.58      0.95      0.72     26043
           1       0.87      0.30      0.44     26043

    accuracy                           0.63     52086
   macro avg       0.72      0.63      0.58     52086
weighted avg       0.72      0.63      0.58     52086
```

```
Results for DecisionTreeClassifier:
              precision    recall  f1-score   support

           0       0.98      0.98      0.98     26043
           1       0.98      0.98      0.98     26043

    accuracy                           0.98     52086
   macro avg       0.98      0.98      0.98     52086
weighted avg       0.98      0.98      0.98     52086


Results for XGBClassifier:
              precision    recall  f1-score   support

           0       0.92      0.91      0.92     26043
           1       0.91      0.92      0.92     26043

    accuracy                           0.92     52086
   macro avg       0.92      0.92      0.92     52086
weighted avg       0.92      0.92      0.92     52086
```

After the evaluation of model accuracy RandomForestClassifier and BaggingClassifier is selected for model hyper parameter tuning using GridSearchCV with different combinations of parameters and the accuracy remains the same.

| | A | Logistic Regression | LinearSVC | Random Forest Classifier | Gradient Boosting Classifier | Bagging Classifier | SVC | Decision Tree Classifier | XGB Classifier |
|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | |
| 2 | 0 | 0.612209637 | 0.603378768 | 0.773085045 | 0.809944327 | 0.776924554 | 0.621616433 | 0.756575158 | 0.809944327 |
| 3 | 1 | 0.621808409 | 0.769437512 | 0.782491841 | 0.816471492 | 0.782875792 | 0.629871376 | 0.7640622 | 0.813783836 |
| 4 | 2 | 0.617584949 | 0.721827606 | 0.77750048 | 0.813399885 | 0.782875792 | 0.624880015 | 0.760798618 | 0.810904204 |
| 5 | 3 | 0.620464581 | 0.690151661 | 0.918410443 | 0.896333269 | 0.91898637 | 0.625647917 | 0.90132463 | 0.919946247 |
| 6 | 4 | 0.628335573 | 0.638702246 | 0.944135151 | 0.91898637 | 0.942599347 | 0.632942983 | 0.926473411 | 0.946054905 |
| 7 | 5 | 0.623152237 | 0.810712229 | 0.941447495 | 0.908619697 | 0.938759839 | 0.631791131 | 0.919562296 | 0.942215396 |
| 8 | 6 | 0.611367127 | 0.810291859 | 0.934907834 | 0.90437788 | 0.934331797 | 0.618471582 | 0.909178187 | 0.939900154 |
| 9 | 7 | 0.6140553 | 0.783218126 | 0.944892473 | 0.911290323 | 0.93874808 | 0.621543779 | 0.919354839 | 0.941436252 |
| 10 | 8 | 0.625576037 | 0.78859447 | 0.946620584 | 0.902841782 | 0.942780338 | 0.631528418 | 0.916282642 | 0.935675883 |
| 11 | 9 | 0.615399386 | 0.810099846 | 0.950076805 | 0.913402458 | 0.947580645 | 0.623463902 | 0.926459293 | 0.947004608 |

**Table: 10 Fold Cross Validation**