In this project, you are going to work on the The "Census Income" data set from the UCI Machine Learning Repository that contains the income information for over 48,000 individuals taken from the 1994 US census. For more details about this dataset, you can refer to the following link: https://archive.ics.uci.edu/ml/datasets/census+income

Problem Statement:

In this project, initially you need to preprocess the data and then develop an understanding of different features of the data by performing exploratory analysis and creating visualizations. Further, after having sufficient knowledge about the attributes you will perform a predictive task of classification to predict whether an individual makes over 50K a year or less, by using different Machine Learning Algorithms.

1. Data Preprocessing:

fnlwgt education education-num marital-status occupation relationship

a) Replace all the missing values with NA.

```
b) Remove all the rows that contain NA values.
           import pandas as pd
           import numpy as np
import matplotlib.pyplot as plt
          #Read the Census Income dataset
df = pd.read_csv(r'C:\Jagan\Personal\DS&AI Certification\Census_Income_Project\census-income.csv')
 In [3]: df.head()
                       workclass fnlwgt education education-num marital-status
           age
                                                                                     occupation relationship race sex capital-qain capital-loss hours-per-week native-country
                                                          13 Never-married
                                                                                                                                  2174
          0 39
                       State-gov 77516 Bachelors
                                                                                      Adm-clerical Not-in-family White Male
                                                                                                                                                0
                                                                                                                                                              40
                                                                                                                                                                    United-States <=50K
          1 50 Self-emp-not-inc 83311 Bachelors
                                                                                                                                  0
                                                                                                                                                           13 United-States <=50K
                                                         13 Married-civ-spouse Exec-managerial
                                                                                                  Husband White Male
                                                    9
                                                                                                                                    0
                                                                                                                                                0
          2 38
                        Private 215646 HS-grad
                                                                      Divorced Handlers-cleaners Not-in-family White Male
                                                                                                                                                              40 United-States <=50K
                    Private 234721
                                         11th
                                                                                                                                                          40 United-States <=50K
          3 53
                                                    7 Married-civ-spouse Handlers-cleaners Husband Black Male
          4 28
                        Private 338409 Bachelors
                                                            13 Married-civ-spouse
                                                                                     Prof-specialty
                                                                                                        Wife Black Female
                                                                                                                                                              40
                                                                                                                                                                          Cuba <=50K
 In [4]: df.describe()
                                   fnlwgt education-num capital-gain capital-loss hours-per-week
          count 32561.000000 3.256100e+04 32561.000000 32561.000000 32561.000000
                                                                                   32561.000000
                 38.581647 1.897784e+05 10.080679 1077.648844 87.303830
                                                                                   40.437456
             std
                   13.640433 1.055500e+05
                                                2.572720 7385.292085 402.960219
                                                                                       12 347429
            min
                 17.000000 1.228500e+04
                                           1.000000 0.000000 0.000000
                                                                                     1.000000
           25%
                   28.000000 1.178270e+05
                                                9.000000
                                                            0.000000
                                                                         0.000000
                                                                                       40.000000
           50%
                 37.000000 1.783560e+05
                                               10.000000 0.000000
                                                                         0.000000
                                                                                       40.000000
           75%
                 48.000000 2.370510e+05
                                               12.000000
                                                            0.000000
                                                                         0.000000
                                                                                       45 000000
                 90.000000 1.484705e+06
           max
                                              16.000000 99999.000000 4356.000000
                                                                                       99.000000
 In [5]: df.columns
In [6]: df.info()
          <class 'pandas.core.frame.DataFrame':
          RangeIndex: 32561 entries, 0 to 32560
Data columns (total 15 columns):
# Column Non-Null Count Dtype
              age
workclass
fnlwgt
education
                                 32561 non-null int64
                                 32561 non-null object
                32561 non-null
                                                  object
                capital-gain
                                 32561 non-null
          10 capital-gain 3/2561 non-null int64
11 capital-loss 3/2561 non-null int64
12 hours-per-week 3/2561 non-null int64
13 native-country 3/2561 non-null object
14 3/2561 non-null object
dtypes: int64(6), object(9)
memory usage: 3.7+ MB
In [14]: df.isna().sum()
Out[14]: age workclass
           fnlwgt
education
education-num
marital-status
           occupation
           relationship
           capital-gain
capital-loss
hours-per-week
native-country
          dtype: int64
          df.isnull().sum()
Out[15]: age workclass
```

```
sex
       capital-gain
       capital-loss
      hours-per-week
native-country
      dtype: int64
      df_columns = df.columns
      df_columns
' '],
dtype='object')
In [138-
#replacing space before the column names
df.columns = df.columns.str.replace(' ', '')
In [21]: df.columns
In [23]: df.shape
Out[23]: (32561, 15)
In [42]: df.columns.values[14] = "YearlyIncome"
In [43]: df.columns
dtype='object'
```

2. Data Manipulation:

- a) Extract the "education" column and store it in "census_ed".
- b) Extract all the columns from "age" to "relationship" and store it in "census_seq".
- c) Extract the column number "5", "8", "11" and store it in "census_col".
- d) Extract all the male employees who work in state-gov and store it in "male_gov".
- e) Extract all the 39 year olds who either have a bachelor's degree or who are native of the United States and store the result in "census_us".
- f) Extract 200 random rows from the "census" data frame and store it in "census_200".
- g) Get the count of different levels of the "workclass" column.
- h) Calculate the mean of the "capital.gain" column grouped according to "workclass".
- i) Create a separate dataframe with the details of males and females from the census data that has income more than 50,000.
- j) Calculate the percentage of people from the United States who are private employees and earn less than 50,000 annually.
- k) Calculate the percentage of married people in the census data.
- I) Calculate the percentage of high school graduates earning more than 50,000 annually.

```
In [44]: # Extract the "education" column and store it in "census_ed" census_ed = df.loc[:,('education')]

In [45]: type(census_ed)

Out[45]: pandas.core.series.Series

In [46]: census_ed.shape

Out[46]: (32561,)
```

b) Extract all the columns from "age" to "relationship" and store it in "census_seq"

Out[48]: (32561, 8)

c) Extract the column number "5", "8", "11" and store it in "census_col"

```
In [51]: census_col = df.iloc[:,([5,8,11])]
In [53]: census_col.columns
Out[53]: Index(['marital-status', 'race', 'capital-loss'], dtype='object')
In [54]: census_col.shape
Out[54]: (32561, 3)
```

```
d) Extract all the male employees who work in state-gov and store it in "male_gov".
In [58]:
Out[58]: array([' Male', ' Female'], dtype=object)
        df['workclass'].unique()
Out[59]: array([' State-gov', ' Self-emp-not-inc', ' Private', ' Federal-gov', ' Local-gov', ' ?', ' Self-emp-inc', ' Without-pay', ' Never-worked'], dtype=object)
In [66]: male_gov = df.loc[(df['sex'] ==' Male') & (df['workclass'] == ' State-gov'),]
             age workclass fnlwgt education education-num marital-status occupation relationship
                                                                                                race sex capital-gain capital-loss hours-per-week native-country YearlyIncome
                                                13 Never-married Adm-clerical Not-in-family
                                                                                                                     0 40 United-States
           0 39 State-gov 77516
                                 Bachelors
                                                                                               White Male
                                                                                                              2174
                                                                                                                                 40 India
         11 30 State-gov 141297 Bachelors
                                               13 Married-civ-spouse Prof-specialty Husband Asian-Pac-Islander Male
                                                                                                                                                        >50K
          34 22 State-gov 311512 Some-college
                                                10 Married-civ-spouse Other-service
                                                                                Husband
                                                                                              Black Male
                                                                                                                0 0
                                                                                                                                   15 United-States
                                                                                                                                                        <=50K
        48 41 State-gov 101603 Assoc-voc 11 Married-civ-spouse Craft-repair Husband
                                                                                               White Male
                                                                                                                0 0 40 United-States
                                                                                                                                                        <=50K
         123 29 State-gov 267989 Bachelors 13 Married-civ-spouse Prof-specialty
                                                                                Husband
                                                                                               White Male
                                                                                                                                   50 United-States
                                                                                                                                                        >50K
        32163 36 State-gov 135874 Bachelors 13 Married-civ-spouse Sales Husband
                                                                                               White Male
                                                                                                                0 0
                                                                                                                                    40 United-States
                                                                                                                                                        <=50K
        32241 45 State-gov 231013 Bachelors 13 Divorced Protective-serv Not-in-family
                                                                                                                0 0 40 United-States
                                                                                            White Male
                                                                                                                                                        <=50K
        32321 54 State-gov 138852 HS-grad 9 Married-civ-spouse Prof-specialty Husband
                                                                                                                0 0
                                                                                              White Male
                                                                                                                                   40 United-States
                                                                                                                                                        <=50K
                                               10 Divorced Adm-clerical Own-child
                                                                                                                0 0 40 United-States
        32324 42 State-gov 138162 Some-college
                                                                                              White Male
                                                                                                                                                        <=50K
                                 HS-grad 9 Married-civ-spouse Craft-repair Husband
        32360 58 State-gov 200316
                                                                                               White Male
                                                                                                                                  40 United-States
                                                                                                                                                        <=50K
       809 rows × 15 columns
In [67]: male_gov.shape
Out[67]: (809, 15)
```

e) Extract all the 39 year olds who either have a bachelor's degree or who are native of the United States and store the result in "census_us".

```
Dut (18] | df (18ge) | werkclass', 'fmlugt', 'edication', 'edication',
```

White Male

9 Divorced Exec-managerial Not-in-family White Male 0 0 80 United-States <=50K

2174

40 United-States

13

Never-married

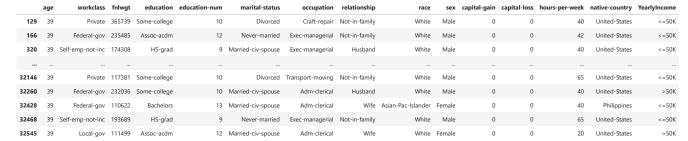
Adm-clerical Not-in-family

0 39

State-gov 77516

28 39 Private 367260 HS-grad

Bachelors



759 rows × 15 columns

f) Extract 200 random rows from the "census" data frame and store it in "census_200".

```
df.shape
Out[83]: (32561, 15)
        def Rand(start, end, num):
            for j in range(num):
    res.append(random.randint(start, end))
           return res
         census_200 = df.iloc[(Rand(0, 32561, 200))]
        census 200
              age workclass fnlwgt education education-num marital-status
                                                                      occupation relationship race
                                                                                                   sex capital-gain capital-loss hours-per-week native-country YearlyIncome
Out[88]:
                                           10 Never-married
                                                                                                                   0
                                                                                                                                  49 United-States
        21330 33
                   Private 178449 Some-college
                                                                       Adm-clerical Not-in-family White Male
                                                                                                              0
        32476 35 Private 30673 12th
                                                8 Married-civ-spouse
                                                                     Craft-repair Husband White
                                                                                                                                 84 United-States
        13135 24 State-gov 275421 Some-college
                                                        Never-married Machine-op-inspct Own-child White Female
                                                                                                                                   40 United-States
        21957 27 State-gov 346406 Bachelors 13 Never-married Prof-specialty Unmarried White
                                                                                                           0 0 50 United-States
                                     10th
         4373 35
                   Private 111128
                                                  6 Never-married
                                                                        Adm-clerical Own-child White
                                                                                                                                   40 United-States
                                                                                                                                                       <=50K
                                                  11 Married-civ-spouse
        25065 37
                   Private 249208
                                Assoc-voc
                                                                      Prof-specialty Husband White Male
                                                                                                                                   48 United-States
        9177 36 Private 114605 Assoc-voc 11 Married-civ-spouse Adm-clerical Husband White
                                                                                                   Male 0 0 40 United-States
                                                                                                                                                        >50K
                                                   13 Married-civ-spouse Exec-managerial Husband White
        16829 44 State-gov 175696
                                                                                                                                   40 United-States
                                  Bachelors
        17795 67 Private 172756 1st-4th
                                                2 Widowed Machine-op-inspct Not-in-family White Female
                                                                                                             2062 0 34 Ecuador
                                                                                                                                                        <=50K
        24079 34 Private 115858
                                  HS-grad
                                                          Divorced
                                                                       Adm-clerical Own-child White Female
                                                                                                             0 0 40 United-States
       200 rows × 15 columns
        census 200 shape
Out[89]: (200, 15)
        census_200.head()
             age workclass fnlwgt education education-num marital-status
                                                                      occupation relationship race
                                                                                                   sex capital-gain capital-loss hours-per-week native-country YearlyIncome
        21330 33 Private 178449 Some-college
                                           10 Never-married
                                                                        Adm-clerical Not-in-family White
                                                 8 Married-civ-spouse
                                                                      Craft-repair Husband White
             24 State-gov 275421 Some-college
                                                         Never-married Machine-op-inspct
        21957 27 State-gov 346406 Bachelors 13 Never-married Prof-specialty Unmarried White
```

g) Get the count of different levels of the "workclass" column.

h) Calculate the mean of the "capital.gain" column grouped according to "workclass".

Federal-gov 833.23292
Local-gov 880.202580
Never-worked 0.000000
Private 889.217792
Self-emp-inc 4875.693548
Self-emp-not-inc 1886.061787
State-gov Without-pay 487.8571436
Name: capital-gain, dtype: float64

i) Create a separate dataframe with the details of males and females from the census data that has income more than 50,000.

```
columns = df.columns
         columns
Out[3]: Index(['age', 'workclass', 'fnlwgt', 'education', 'education-num', 'marital-status', 'occupation', 'relationship', 'race', 'sex', 'capital-gain', 'capital-loss', 'hours-per-week', 'native-country',
              dtype='object')
In [180...
          #replacing space before the column names
         df.columns = df.columns.str.replace(' ', '')
In [9]: df.columns.values[14] = "YearlyIncome"
In [14]: df.columns
dtype='object')
         df.rename(columns={df.columns[14]: "YearlyIncome"},inplace=True)
print(df.columns)
        dtype='object')
        df.sex.unique()
Out[8]: array([' Male', ' Female'], dtype=object)
         df['YearlyIncome'].value_counts()
         <=50K 24720
        Name: YearlyIncome, dtype: int64
         males 50k = df.loc[(df['sex'] == ' Male') & (df['YearlvIncome'] == ' >50K').]
In [50]:
         males_50k.head()
                                                              marital-status occupation relationship
           age workclass fnlwgt education education-num
                                                                                                         race sex capital-gain capital-loss hours-per-week native-country YearlyIncome
         7 52 Self-emp-not-inc 209642
                                      HS-grad
                                                       9 Married-civ-spouse Exec-managerial
                                                                                                           White Male
                                                                                                                             0
                                                                                                                                                     45 United-States
                                                                                           Husband
        9 42 Private 159449 Bachelors
                                                     13 Married-civ-spouse Exec-managerial Husband
                                                                                                           White Male 5178 0 40 United-States
                     Private 280464 Some-college
                                                                                                          Black Male
                                                       10 Married-civ-spouse Exec-managerial
                                                     13 Married-civ-spouse Prof-specialty Husband Asian-Pac-Islander Male 0 0 40 India >50K
         11 30 State-gov 141297 Bachelors
                                                       11 Married-civ-spouse Craft-repair
                     Private 121772 Assoc-voc
        females_50k = df.loc[(df['sex'] == ' Female') & (df['YearlyIncome'] == ' >50K'),]
In [52]: females_50k.head()
                     workclass fnlwgt education education-num marital-status occupation relationship race sex capital-gain capital-loss hours-per-week native-country YearlyIncome
                     Private 45781 Masters
                                                             Never-married Prof-specialty Not-in-family White Female
                                                    14 Divorced Exec-managerial Unmarried White Female 0 0 45 United-States >50K
        19 43 Self-emp-not-inc 292175 Masters
        67 53 Private 169846 HS-grad 9 Married-civ-spouse Adm-clerical Wife White Female 0 0 40 United-States
                      Private 343591
                                    HS-grad
                                                                            Craft-repair Not-in-family White Female
         print("Count of records where sex = male and Income >=50k : ", males_50k.shape[0])
print("Count of records where sex = female and Income >=50k : ", females_50k.shape[0])
        Count of records where sex = male and Income >=50k : 6662 Count of records where sex = female and Income >=50k : 1179
```

j) Calculate the percentage of people from the United States who are private employees and earn less than 50,000 annually.

```
'Trinadad&Tobago', 'Greece', 'Nicaragua', 'Vietnam', 'Hor
'Ireland', 'Hungary', 'Holand-Netherlands'], dtype=object)
In [45]: df.workclass.unique()
Out[45]: array(['State-gov', 'Self-emp-not-inc', 'Private', 'Federal-gov', 'Local-gov', '?', 'Self-emp-inc', 'Without-pay', 'Never-worked'], dtype=object)
In [47]: df['YearlyIncome'].unique()
Out[47]: array([' <=50K', ' >50K'], dtype=object)
          df_us_under_50k_private = df.loc[(df['native-country'] == ' United-States') & (df['YearlyIncome'] == ' <=50K') & (df['workclass'] == ' Private'),]
df_us_under_50k_private.head()</pre>
Out[55]: age workclass fnlwgt education education-num marital-status occupation relationship race sex capital-gain capital-loss hours-per-week native-country YearlyIncome
                                                       9 Divorced Handlers-cleaners Not-in-family White Male
          2 38 Private 215646 HS-grad
                                                                                                                                                     40 United-States
          3 53 Private 234721 11th
                                                                                                                           0 0 40 United-States
                                                      7 Married-civ-spouse Handlers-cleaners Husband Black Male
                                                     14 Married-civ-spouse Exec-managerial Wife White Female
          5 37 Private 284582 Masters
                                                                                                                                                    40 United-States
                                                   13 Never-married Adm-clerical Own-child White Female 0 0 30 United-States
          12 23 Private 122272 Bachelors
          13 32 Private 205019 Assoc-acdm 12 Never-married
                                                                                   Sales Not-in-family Black Male 0 0 50 United-States
          percentage_us_under_50k_private = (len(df_us_under_50k_private) * 100 )/ len(df)
print('percentage of US, under_50k income and private employed: ', percentage_us_under_50k_private)
          percentage of US, under 50k income and private employed: 47,891649519363654
```

k) Calculate the percentage of married people in the census data.

```
In [59]: df['marital-status'].unique()
In [60]: df['marital-status'].value_counts()
Out[60]: Married-civ-spouse
      Widowed
Married-spouse-absent
                     993
418
      Married-AF-spouse
      Name: marital-status, dtype: int64
      Percentage of married people : 47.34805442093302
```

I) Calculate the percentage of high school graduates earning more than 50,000 annually.

```
In [71]: df.education.unique()
Out[71]: array([' Bachelors', ' HS-grad', ' 11th', ' Masters', ' 9th', ' Some-college', ' Assoc-acdm', ' Assoc-voc', ' 7th-8th', ' Doctorate', ' Prof-school', ' 5th-6th', ' 10th', ' 1st-4th', ' Preschool', ' 12th'], dtype=object)
In [72]: df.education.value_counts()
Out[72]: HS-grad
                  Some-college
                  Bachelors
                  Masters
                 Masters
Assoc-voc
11th
Assoc-acdm
10th
7th-8th
                  Prof-school
                  12th
                  Doctorate
5th-6th
1st-4th
                  Preschool
                                                  51
                Name: education, dtype: int64
                hs_grad_income_gt_50k = len (df.loc[(df['education'] == ' HS-grad') & (df['YearlyIncome'] == ' <=50K'),])
print("Highschool graduates with income >50k : ", hs_grad_income_gt_50k)
print("Total records count : ", len(df))
                 print("percentaage of Highschool graduates with income >50k : " ,(hs_grad_income_gt_50k *100 )/ len(df))
               Highschool graduates with income >50k : 8826
Total records count : 32561
percentaage of Highschool graduates with income >50k : 27.106047111575197
```

3. Linear Regression:

- a) Build a simple linear regression model as follows:
- Divide the dataset into training and test sets in 70:30 ratio.
- Build a linear model on the test set where the dependent variable is "hours.per.week" and the independent variable is "education.num".
- Predict the values on the train set and find the error in prediction.
- Find the root-mean-square error (RMSE).

```
In [77]: df.columns
In [87]:     x = df[['education-num']]
     y = df[['hours-per-week']]
In [88]: type(x), type(y)
Out[88]: (pandas.core.frame.DataFrame, pandas.core.frame.DataFrame)
In [80]:
            from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
In [81]: | lr=LinearRegression()
In [89]:
           # Divide the dataset into training and test sets in 70:30 ratio.
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.30,random_state=1)
In [90]:
            #fit the model on training data
            lr.fit(x_train,y_train)
Out[90]: LinearRegression()
            #Predict the values on the train set
            y_pred=lr.predict(x_test)
In [92]: #find the error in prediction
           error = y_test - y_pred
error
                 hours-per-week
            9646
           709
            7385
           16671
           21932
           29663
                        -1.832763
           29310
           19491
                        -1.098059
            2861
                         5.514277
          9769 rows × 1 columns
           print('mean_squared_error :',mean_squared_error(y_test,y_pred))
            print('root-mean-square error :',np.sqrt(mean_squared_error(y_test,y_pred)))
```

4. Logistic Regression:

mean_squared_error : 147.15261838664162 root-mean-square error : 12.130647896408568

- a) Build a simple logistic regression model as follows:
- Divide the dataset into training and test sets in 65:35 ratio.
- Build a logistic regression model where the dependent variable is "X" (yearly income) and the independent variable is "occupation".
- Predict the values on the test set.
- Build a confusion matrix and find the accuracy.

```
1843
            Transport-moving
Handlers-cleaners
Farming-fishing
Tech-support
Protective-serv
            Priv-house-serv
             Armed-Forces
           dtype: int64
In [97]: #occupation is indpendent
            x=df['occupation'].replace('?','Prof-specialty')
x=pd.DataFrame(x)
In [99]: df['YearlyIncome'].unique()
Out[99]: array([' <=50K', ' >50K'], dtype=object)
In [100... #income is dependent
            y=df['YearlyIncome'].replace(' <=50K',0).replace(' >50K',1)
           y=pd.DataFrame(y)
In [101... x.head()
Out[101...
In [102... y.head()
Out[102...
           YearlyIncome
          1
In [105... le=LabelEncoder()
           x=le.fit_transform(x)
           C:\Users\577346744\Anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for
           example using ravel().
return f(*args, **kwargs)
In [106... type(x)
Out[106... numpy.ndarray
In [109...
          x=nd.DataFrame(x)
In [110... type(x)
Out[110... pandas.core.frame.DataFrame
             x\_train, x\_test, y\_train, y\_test=train\_test\_split(x,y,test\_size=.35, random\_state=1) \\ 1o=LogisticRegression() 
            lo.fit(x train, v train)
            y_pred=lo.predict(x_test)
            print('confusion_matrix : ')
           print(confusion_matrix(y_pred,y_test))
print('accuracy_score : ',accuracy_score(y_test,y_pred))
           confusion_matrix :
[[8800 2597]
[ 0 0]]
           accuracy_score :
                               0.7721330174607353
           C:\Users\S77346744\Anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for
```

- 4 b)Build a multiple logistic regression model as follows:
- Divide the dataset into training and test sets in 80:20 ratio.
- Build a logistic regression model where the dependent variable is "X" (yearly income) and independent variables are "age", "workclass", and "education".
- Predict the values on the test set.

example using ravel().
 return f(*args, **kwargs)

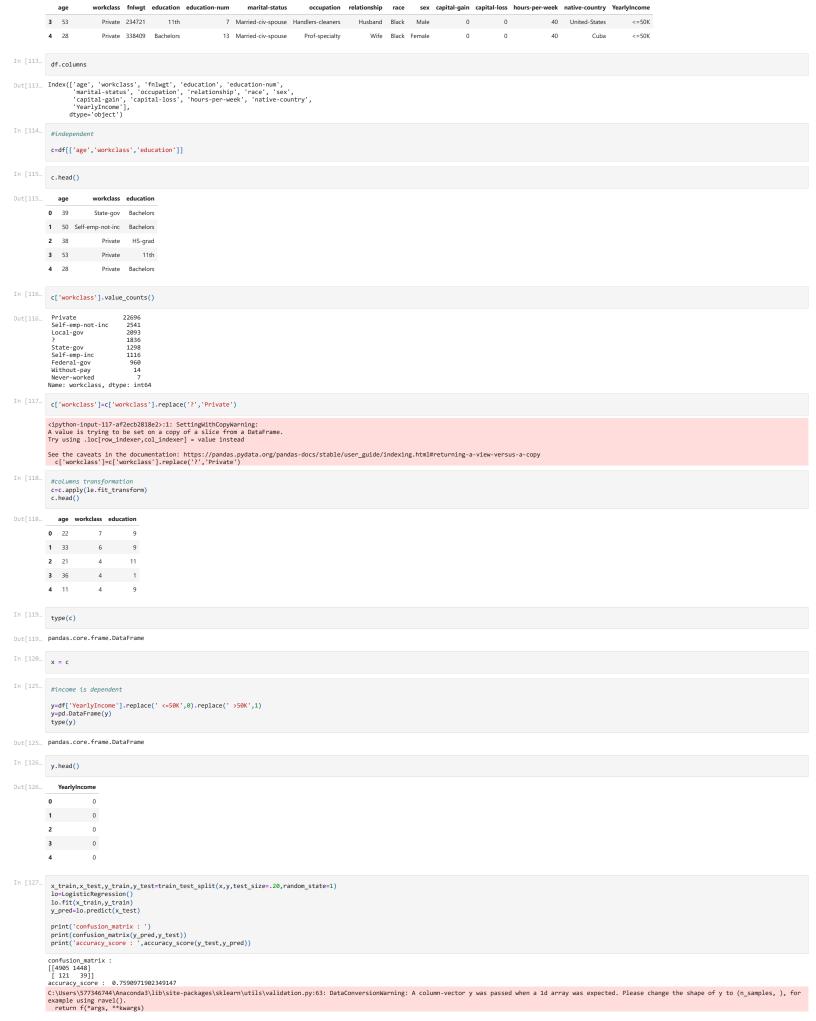
Sales Other-service

Machine-op-inspct

3295

• Build a confusion matrix and find the accuracy.

In [112	df	head()														
Out[112		ige	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	hours-per-week	native-country	YearlyIncome
	0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	40	United-States	<=50K
	1	50 Self-	-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<=50K
	2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<=50K



5. Decision Tree:

- a) Build a decision tree model as follows:
- Divide the dataset into training and test sets in 70:30 ratio.
- Build a decision tree model where the dependent variable is "X" (Yearly Income) and the rest of the variables as independent variables.
- Predict the values on the test set.

In [188... df['marital-status'].unique()

dtype='object')

Out[191... array([' <=50K', ' >50K'], dtype=object)

In [191... df.YearlyIncome.unique()

In [189... df.columns

Out[194...

age

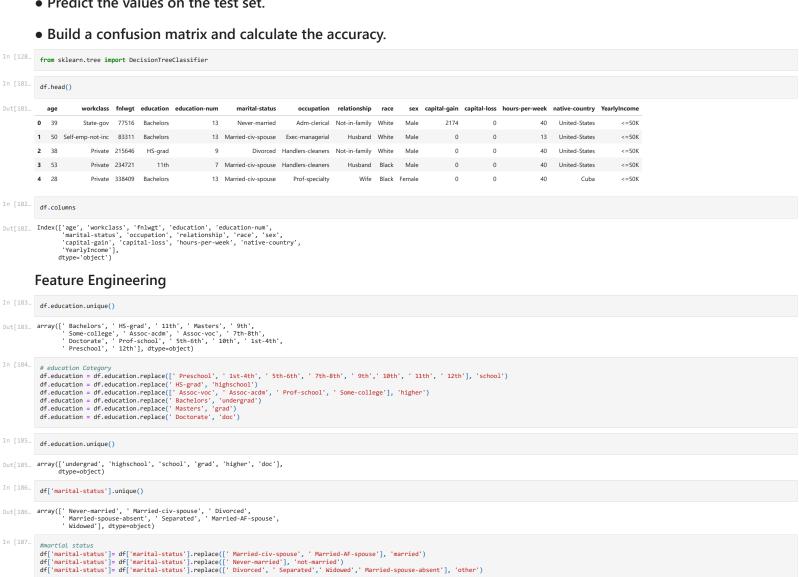
Out[188__array(['not-married', 'married', 'other'], dtype=object)

#"repetcing'
df['workclass']-df['workclass'].replace('?','Private')
df['occupation']=df['occupation'].replace('?','Prof-specialty')
df['native-country']=df['native-country'].replace('?','United-States')

workclass fnlwgt education education-num marital-status

In [190— # Rename single column by Index df.rename(columns={df.columns[14]: "YearlyIncome"},inplace=True)

df.income = df.YearlyIncome.replace(' <=50K', 0
df.income = df.YearlyIncome.replace(' >50K', 1)



<ipython-input-192-d015b1d7e588>:2: UserWarning: Pandas doesn't allow columns to be created via a new attribute name - see https://pandas.pydata.org/pandas-docs/stable/indexing.html#attribute-access
df.income = df.YearlyIncome.replace(' <=50K', 0)</pre>

occupation relationship race sex capital-gain capital-loss hours-per-week native-country YearlyIncome

•	-	33	-	tate-	gov	,,,,,,,,	dilu	uergrau		13 1100	t-illalilleu	Adili-cierica	at IN	101-111-12	arrinty	vviiite	iviale	4	.174	J	-	+0	Officeu-States	V=30
1		50 S	elf-emp	-not-	inc	83311	unc	dergrad		13	married	Exec-manageria	al	Hush	band	White	Male		0	0	1	13	United-States	<=50K
2	3	38		Priv	ate 2	215646	high	hschool		9	other	Handlers-cleaner	rs N	ot-in-fa	amily	White	Male		0	0	4	40	United-States	<=50K
3	į	53		Priv	ate 2	234721		school		7	married	Handlers-cleaner	rs	Husl	band	Black	Male		0	0	4	40	United-States	<=50K
4	2	28		Priv	ate 3	338409	und	dergrad		13	married	Prof-specialt	ty		Wife	Black	Female		0	0	4	40	Cuba	<=50K
5 d-	f.c	colum	ıns																					
5 In	de	'r	narita	l-st l-ga Inco	atus in', me']	', 'od 'cap:	ccupa	/gt', 'educa htion', 'rel ·loss', 'hou	latio	nship', 'ra	ace', 's	sex',												
d-	ack f=c	kup=c		y()		ansfor	m)																	
97	aç	ge w	orkcla	ss f	nlwgt	educ	ation	education-r	num	marital-statu	ıs occup	oation relations	hip	race :	sex (capital-g	jain car	oital-loss	hours-per-w	eek n	native-country	y Ye	earlyIncome	
0	2	22		7	2671		5		12		1	1	1	4	1		25	0		39	39	9	0	
1	3	33		6	2926		5		12		0	4	0	4	1		0	0		12	39	9	0	
2	2	21		4	14086		3		8		2	6	1	4	1		0	0		39	39	9	0	
3	3	36		4	15336		4		6		0	6	0	2	1		0	0		39	39	9	0	
4	1	11		4	19355		5		12		0	10	5	2	0		0	0		39	5	5	0	
X:	=d1	f.ilo pende	ndent oc[:,: ent oc[:,-																					
99 t	уре	e(x)																						
9 pa	nda	as.co	ore.fr	ame.	Data	Frame																		
d· d· y.	t=C t.f _pr _rir	Decis fit() red=c nt('c	ionTr _trai t.pre confus onfusi	eeCl n,y_ dict ion_ on_m	assintrain (x_te matri	fier() n) est) ix : ' x(y_pr	') red,y	rain_test_s v_test)) vacy_score(y				∂,random_state	:=1)											
[[[ac	654 100 cui	47 8 03 13 racy_	89]] _score	::	0.81	236564 Ore																		

occupation relationship race sex capital-gain capital-loss hours-per-week native-country YearlyIncome

0

40 United-States

2174

Adm-clerical Not-in-family White Male

age

0 39

a) Build a random forest model as follows:

workclass fnlwgt education education-num marital-status

13 not-married

State-gov 77516 undergrad

- Divide the dataset into training and test sets in 80:20 ratio.
- Build a random forest model where the dependent variable is "X" (Yearly Income) and the rest of the variables as independent variables and number of trees as 300.
- Predict values on the test set
- Build a confusion matrix and calculate the accuracy

```
In [201... from sklearn.ensemble import RandomForestClassifier
In [202- x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.20,random_state=1) rf=RandomForestClassifier(n_estimators=300)
             rf.fit(x train,y train)
             y_pred=rf.predict(x_test)
             print('confusion matrix : ')
             print(confusion_matrix(y_pred,y_test))
print('accuracy_score : ',accuracy_score(y_test,y_pred))
            confusion_matrix :
[[4658 523]
[ 368 964]]
```

- 7. For this problem, use the population dataset, and perform the following:
- 1. EDA on the time series to find trends and seasonality.
- 2. Forecast the population on the given dataset for the next 6 months.

This problem was solved in Colab as prophet not working on Jupyter. Shared another file for that

```
popdata = pd.read_csv(r'C:\Jagan\Personal\DS&AI Certification\Census_Income_Project\popdata.csv')
In [204...
         popdata.head()
```

```
value
                                                  date
                  0 127299.0 1952-01-01
                  1 127517.0 1952-02-01
                   2 127721.0 1952-03-01
                  3 127933.0 1952-04-01
                   4 128130.0 1952-05-01
In [205... popdata.columns
Out[205... Index(['value', 'date'], dtype='object')
 In [206...
                   popdata.describe()
 Out[206...
                                             value
                                   816.000000
                   mean 214837.767826
                       std 50519.140567
                      min 127299.000000
                     25% 172715.250000
                     50% 210547.500000
                     75% 260354.250000
                      max 301299.946000
 In [207...
                   popdata.info()
                  <class 'pandas.core.frame.DataFrame'>
RangeIndex: 816 entries, 0 to 815
Data columns (total 2 columns):
# Column Non-Null Count Dtype
                  0 value 816 non-null float64
1 date 816 non-null object
dtypes: float64(1), object(1)
memory usage: 12.9+ KB
In [209...
                   popdata.isnull().sum()
Out[209... value 0
                  dtvpe: int64
 In [210...
                   from datetime import datetime
                   popdata['Value']=popdata['value']
popdata=popdata.drop(columns=['value'])
popdata.head()
                                 date Value
                  0 1952-01-01 127299.0
                  1 1952-02-01 127517.0
                  2 1952-03-01 127721.0
                  3 1952-04-01 127933.0
                  4 1952-05-01 128130.0
In [212... popdata.columns=['ds','y']
In [213... popdata.head()
                                    ds
                                                      у
                  0 1952-01-01 127299.0
                  1 1952-02-01 127517.0
                   2 1952-03-01 127721.0
                  3 1952-04-01 127933.0
                   4 1952-05-01 128130.0
In [227... from fbprophet import Prophet
                   ModuleNotFoundError Traceback (most recent call last) <ipython-input-227-f503e9c6cf11> in <module>
                     ---> 1 from fbprophet import Prophet
                  ModuleNotFoundError: No module named 'fbprophet'
In [226...
                   pip install --upgrade setuptools
                    Requirement already satisfied: setuptools in c:\users\577346744\anaconda3\lib\site-packages (52.0.0.post20210125)
                 Requirement already satisfied: setuptools in c:\users\577346744\anaconda3\lib\site-packages (52.0.0.post20210125)
Collecting setuptools
Downloading setuptools-62.3.2-py3-none-any.whl (1.2 MB)
Installing collected packages: setuptools
Attempting uninstall: setuptools
Found existing installation: setuptools 52.0.0.post20210125
Uninstalling setuptools-52.0.0.post20210125:
Successfully uninstalled setuptools-52.0.0.post20210125
Successfully uninstalled setuptools-52.0.0.post20210125
Successfully installed setuptools-62.3.2
Note: you may need to restart the kernel to use updated packages.
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency conflicts. spyder 4.2.5 requires pyqtso.13.y which is not installed.
conda-repo-cli 1.0.4 requires pathlib, which is not installed.
                   pip install fbprophet
                       Using cached fbprophet-0.7.1.tar.gz (64 kB)
                  Using cached fbprophet-0.7.1.tan.gz (64 kB) Requirement already satisfied: Cython>=0.22 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (0.29.23) Requirement already satisfied: cmdstanpy==0.9.5 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (0.9.5) Requirement already satisfied: pystan>=2.14 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (2.19.1.1) Requirement already satisfied: numpy>=1.15.4 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (1.20.1) Requirement already satisfied: pandas>=1.0.4 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (1.2.4)
```

```
Requirement already satisfied: matplotlib>=2.0.0 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (3.3.4)
Requirement already satisfied: convertdate>=2.1.2 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (2.4.0)
Requirement already satisfied: convertdate>=2.1.2 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (2.4.0)
Requirement already satisfied: holidays>=0.10.2 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (0.13)
Requirement already satisfied: setuptools-git>=1.2 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (1.2)
Requirement already satisfied: python-dateutil>=2.8.0 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (2.5.1)
Requirement already satisfied: pymeus<=1,>=0.3.13 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (4.5.9)
Requirement already satisfied: pymeus<=1,>=0.3.13 in c:\users\577346744\anaconda3\lib\site-packages (from fbprophet) (2.2.4)
Requirement already satisfied: pymeus<=1,>=0.3.13 in c:\users\577346744\anaconda3\lib\site-packages (from holidays>=0.10.2-7bprophet) (0.5.11)
Requirement already satisfied: pymeus<=1,>=0.3.13 in c:\users\577346744\anaconda3\lib\site-packages (from holidays>=0.10.2-7bprophet) (0.2.1)
Requirement already satisfied: pymeus<=1,>=0.3.13 in c:\users\577346744\anaconda3\lib\site-packages (from holidays>=0.10.2-7bprophet) (0.2.1)
Requirement already satisfied: pymeus<=1,>=0.3.13 in c:\users\577346744\anaconda3\lib\site-packages (from marCalendar>=0.9->fbprophet) (0.2.1)
Requirement already satisfied: pymeus<=1,>=0.3.13 in c:\users\577346744\anaconda3\lib\site-packages (from matplotlib>=2.0.9-7bprophet) (0.2.1)
Requirement already satisfied: pymeus<=1,>=0.3.13 in c:\users\577346744\anaconda3\lib\site-packages (from matplotlib>=2.0.9-7bprophet) (0.3.1)
Requirement already satisfied: pymeus<=1,>=0.3.13 in c:\users\577346744\anaconda3\lib\site-packages (from matplotlib>=2.0.9-7bprophet) (0.3.1)
Requirement already satisfied: pymeus<=1,>=0.4,1
                                            running build_py
creating build\lib
creating build\lib
creating build\lib\fbprophet
creating build\lib\fbprophet\creating build\lib\fbprophet\creating build\lib\fbprophet\stan_model
Importing plotly failed. Interactive plots will not work.
INFO:pystan:COMPILING THE C++ CODE FOR MODEL anon_model_f5236004a3fd5b8429270d00efcc0cf9 NOW.
WARNING:pystan:MSVC compiler is not supported
error: Microsoft Visual C++ 14.0 or greater is required. Get it with "Microsoft C++ Build Tools": https://visualstudio.microsoft.com/visual-cpp-build-tools/
                                             ERROR: Failed building wheel for fbprophet
ERROR: Command errored out with exit status 1:

command: 'C:\Users\$77346744\AppData\\\\local\\Temp\pip-install-udhyfcac\\fbprophet_1bf20ddd40624185b7f6634c876e0b28\\\setup.py'''''; f=igetatr(tokenize, '"''open'\"'', '"'''); f-iose(); exec(compile(code, __file__, '"''exec''''))' install --record 'C:\Users\$77346744\AppData\\\local\\Temp\pip-record-vkrc35ng\install-record.txt'

--single-version-externally-managed --compile --install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-version-externally-managed --compile --install-badders 'C:\Users\$77346744\AppData\\\local\\Temp\pip-install-udhyfcac\\Temp\pip-install-version-externally-managed --compile --install-badders 'C:\Users\$77346744\AppData\\\Local\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install-udhyfcac\\Temp\pip-install
                                                              cwd: C:\Users\$77346744\
Complete output (11 lines):
running install
running build
running build
reating build\jb
creating build\jb
creating build\jb
creating build\jb\fbprophet\
                                                               Creating build()lib\fbprophet\stan_model
Importing plotly failed. Interactive plots will not work.
INFO:pystan:COMPILING THE C++ CODE FOR MODEL anon_model_f5236004a3fd5b8429270d00efcc0cf9 NOW.
MARNING:pystan:MSVC compiler is not supported
error: Microsoft Visual C++ 14.0 or greater is required. Get it with "Microsoft C++ Build Tools": https://visualstudio.microsoft.com/visual-cpp-build-tools/
                                          ERROR: Command errored out with exit status 1: 'C:\Users\577346744\Anaconda3\python.exe' -u -c 'import sys, setuptools, tokenize; sys.argv[0] = '"""'C:\\Users\577346744\AppData\\Local\\Temp\\pip-install-udhyfca c\\fbprophet_lbf20ddd40624l85b7f6634c876e0b28\\setup.py'""'; _file_='""'C:\\Users\\577346744\\AppData\\Local\\Temp\\pip-install-udhyfcac\\fbprophet_lbf20ddd40624l85b7f6634c876e0b28\\setup.py'""'; _file_='""'C:\\Users\\577346744\\AppData\\Local\\Temp\\pip-install-udhyfcac\\fbprophet_lbf20ddd40624l85b7f6634c876e0b28\\setup.py'""'; _felse('""'\r\n'""); _f.close(); exec(compile(code, _file_, '""'exec'""'))' install --record 'C:\Users\\577346744\AppData\\Local\\Temp\pip-rec ord-vkrc35ng\install-record.txt' --single-version-externally-managed --compile --install-headers 'C:\Users\\577346744\Anacondas\\Throught' Theck the logs for full command output.
In [224...
                                              pip install cryptography
                                             Requirement already satisfied: cryptography in c:\users\577346744\anaconda3\lib\site-packages (3.4.7)
Requirement already satisfied: cffi>=1.12 in c:\users\577346744\anaconda3\lib\site-packages (from cryptography) (1.14.5)
Requirement already satisfied: pycparser in c:\users\577346744\anaconda3\lib\site-packages (from cffi>=1.12->cryptography) (2.20)
Note: you may need to restart the kernel to use updated packages.
                                               model=Prophet()
                                                  model.fit(popdata)
                                                                                                                                                                                                                                              Traceback (most recent call last)
```

NameError: name 'Prophet' is not defined