

In machine learning, Naïve Bayes classification is a straightforward and powerful algorithm for the classification task. In this kernel, I implement Naive Bayes Classification algorithm with Python and Scikit-Learn. I build a Naive Bayes Classifier to predict whether a person makes over 50K a year.

1. Introduction to Naive Bayes algorithm

In machine learning, Naïve Bayes classification is a straightforward and powerful algorithm for the classification task. Naïve Bayes classification is based on applying Bayes’ theorem with strong independence assumption between the features. Naïve Bayes classification produces good results when we use it for textual data analysis such as Natural Language Processing.

Naïve Bayes models are also known as simple Bayes or independent Bayes. All these names refer to the application of Bayes’ theorem in the classifier’s decision rule. Naïve Bayes classifier applies the Bayes’ theorem in practice. This classifier brings the power of Bayes’ theorem to machine learning.

2. Import libraries

```
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt # for data visualization purposes
import seaborn as sns # for statistical data visualization
%matplotlib inline
```

3. Import dataset

```
data = '/content/adult.csv'

df = pd.read_csv(data, header=None, sep=',\s')
```

<ipython-input-2-6bdc593ece69>:3: ParserWarning: Falling back to the 'python' engine because the 'c' engine does not support regex s

df = pd.read_csv(data, header=None, sep=',\s')

4. Exploratory data analysis

```
# view dimensions of dataset

df.shape
```

```
(32561, 15)
```

We can see that there are 32561 instances and 15 attributes in the data set.

```
# preview the dataset

df.head()
```

0

1

2

3

4

5

6

7

8

9

10

11

12

13

14

| | | | | | | | | | | | | | | | |
|---|----|------------------|--------|-----------|----|--------------------|-------------------|---------------|-------|------|------|---|----|---------------|-------|
| 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 |
| 3 | 50 | Private | 224724 | HS-grad | 7 | Married-civ-spouse | Handlers-cleaners | Husband | Black | Male | 0 | 0 | 40 | United-States | <=50K |

Next steps:

Generate code with df

View recommended plots

Rename column names We can see that the dataset does not have proper column names. The columns are merely labelled as 0,1,2.... and so on. We should give proper names to the columns. I will do it as follows:-

```
col_names = ['age', 'workclass', 'fnlwgt', 'education', 'education_num', 'marital_status', 'occupation', 'relationship',
            'race', 'sex', 'capital_gain', 'capital_loss', 'hours_per_week', 'native_country', 'income']
```

```
df.columns = col_names
```

```
df.columns
```

```
Index(['age', 'workclass', 'fnlwgt', 'education', 'education_num',
      'marital_status', 'occupation', 'relationship', 'race', 'sex',
      'capital_gain', 'capital_loss', 'hours_per_week', 'native_country',
      'income'],
      dtype='object')
```

```
# let's again preview the dataset
```

```
df.head()
```

```

age  workclass  fnlwgt  education  education_num  marital_status  occupation  relationship  race  sex  capital_gain  capital_l
0    39   State-gov   77516  Bachelors           13   Never-married  Adm-clerical  Not-in-family  White  Male           2174
1    50  Self-emp-   83311  Bachelors           13   Married-civ-   Exec-        Husband  White  Male              0
      not-inc
2    38   Private  215646  HS-grad            9     Divorced  Handlers-   Not-in-family  White  Male              0
      cleaners

```

Next steps: [Generate code with df](#) [View recommended plots](#)

We can see that the column names are renamed. Now, the columns have meaningful names.

View summary of dataset

```
# view summary of dataset
```

```
df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32561 entries, 0 to 32560
Data columns (total 15 columns):
#   Column                Non-Null Count  Dtype
---  -
0   age                   32561 non-null  int64
1   workclass             32561 non-null  object
2   fnlwgt                32561 non-null  int64
3   education             32561 non-null  object
4   education_num         32561 non-null  int64
5   marital_status        32561 non-null  object
6   occupation            32561 non-null  object
7   relationship          32561 non-null  object
8   race                  32561 non-null  object
9   sex                   32561 non-null  object
10  capital_gain          32561 non-null  int64
11  capital_loss          32561 non-null  int64
12  hours_per_week        32561 non-null  int64
13  native_country        32561 non-null  object
14  income                32561 non-null  object
dtypes: int64(6), object(9)
memory usage: 3.7+ MB

```

We can see that there are no missing values in the dataset

Explore categorical variables

```
# find categorical variables
```

```
categorical = [var for var in df.columns if df[var].dtype=='O']
```

```
print('There are {} categorical variables\n'.format(len(categorical)))
```

```
print('The categorical variables are :\n\n', categorical)
```


```
There are 9 categorical variables
```

```
The categorical variables are :
```



```
['workclass', 'education', 'marital_status', 'occupation', 'relationship', 'race', 'sex', 'native_country', 'income']
```

```
# view the categorical variables
```

```
df[categorical].head()
```



| | workclass | education | marital_status | occupation | relationship | race | sex | native_country | income |
|---|------------------|-----------|--------------------|-------------------|---------------|-------|--------|----------------|--------|
| 0 | State-gov | Bachelors | Never-married | Adm-clerical | Not-in-family | White | Male | United-States | <=50K |
| 1 | Self-emp-not-inc | Bachelors | Married-civ-spouse | Exec-managerial | Husband | White | Male | United-States | <=50K |
| 2 | Private | HS-grad | Divorced | Handlers-cleaners | Not-in-family | White | Male | United-States | <=50K |
| 3 | Private | 11th | Married-civ-spouse | Handlers-cleaners | Husband | Black | Male | United-States | <=50K |
| 4 | Private | Bachelors | Married-civ-spouse | Prof-specialty | Wife | Black | Female | Cuba | <=50K |


Summary of categorical variables There are 9 categorical variables. The categorical variables are given by workclass, education, marital_status, occupation, relationship, race, sex, native_country and income. income is the target variable.

Explore problems within categorical variables

Missing values in categorical variables

```
# check missing values in categorical variables
```

```
df[categorical].isnull().sum()
```



```
workclass      0
education      0
marital_status 0
occupation     0
relationship   0
race           0
sex            0
native_country 0
income         0
dtype: int64
```

We can see that there are no missing values in the categorical variables.

Frequency counts of categorical variables

```
# view frequency counts of values in categorical variables
```

```
for var in categorical:
```

```
    print(df[var].value_counts())
```



```
Taiwan      31
Haiti       44
Iran        43
Portugal    37
Nicaragua   34
Peru        31
France      29
Greece      29
Ecuador     28
Ireland     24
Hong        20
Cambodia    19
Trinidad&Tobago 19
Laos        18
Thailand     18
Yugoslavia  16
Outlying-US(Guam-USVI-etc) 14
Honduras    13
Hungary     13
Scotland    12
Holand-Netherlands 1
Name: count, dtype: int64
income
<=50K      24720
>50K       7841
Name: count, dtype: int64
```

```
# view frequency distribution of categorical variables
```

```
for var in categorical:
```

```
    print(df[var].value_counts()/float(len(df)))
```



```
>>> 0.44081
Name: count, dtype: float64
```

Now, we can see that there are several variables like workclass, occupation and native_country which contain missing values. Generally, the missing values are coded as NaN and python will detect them with the usual command of `df.isnull().sum()`.

Explore workclass variable

```
# check labels in workclass variable

df.workclass.unique()

array(['State-gov', 'Self-emp-not-inc', 'Private', 'Federal-gov',
       'Local-gov', '?', 'Self-emp-inc', 'Without-pay', 'Never-worked'],
      dtype=object)

# check frequency distribution of values in workclass variable

df.workclass.value_counts()

workclass
Private                22696
Self-emp-not-inc       2541
Local-gov              2093
?                      1836
State-gov              1298
Self-emp-inc           1116
Federal-gov            960
Without-pay            14
Never-worked            7
Name: count, dtype: int64
```

We can see that there are 1836 values encoded as ? in workclass variable. I will replace these ? with NaN.

```
# replace '?' values in workclass variable with `NaN`

df['workclass'].replace('?', np.NaN, inplace=True)

# again check the frequency distribution of values in workclass variable

df.workclass.value_counts()

workclass
Private                22696
Self-emp-not-inc       2541
Local-gov              2093
State-gov              1298
Self-emp-inc           1116
Federal-gov            960
Without-pay            14
Never-worked            7
Name: count, dtype: int64
```

Now, we can see that there are no values encoded as ? in the workclass variable.

I will adopt similar approach with occupation and native_country column.

Explore occupation variable

```
# check labels in occupation variable

df.occupation.unique()

array(['Adm-clerical', 'Exec-managerial', 'Handlers-cleaners',
       'Prof-specialty', 'Other-service', 'Sales', 'Craft-repair',
       'Transport-moving', 'Farming-fishing', 'Machine-op-inspct',
       'Tech-support', '?', 'Protective-serv', 'Armed-Forces',
       'Priv-house-serv'], dtype=object)

# check frequency distribution of values in occupation variable

df.occupation.value_counts()

occupation
Prof-specialty         4140
Craft-repair           4099
```

```

Exec-managerial    4066
Adm-clerical       3770
Sales              3650
Other-service      3295
Machine-op-inspct  2002
?                  1843
Transport-moving   1597
Handlers-cleaners  1370
Farming-fishing    994
Tech-support       928
Protective-serv    649
Priv-house-serv    149
Armed-Forces       9
Name: count, dtype: int64

```

We can see that there are 1843 values encoded as ? in occupation variable. I will replace these ? with NaN.

```
# replace '?' values in occupation variable with `NaN`
```

```
df['occupation'].replace('?', np.NaN, inplace=True)
```

```
# again check the frequency distribution of values in occupation variable
```

```
df.occupation.value_counts()
```

```

↔ occupation
Prof-specialty    4140
Craft-repair      4099
Exec-managerial   4066
Adm-clerical      3770
Sales             3650
Other-service     3295
Machine-op-inspct 2002
Transport-moving  1597
Handlers-cleaners 1370
Farming-fishing   994
Tech-support      928
Protective-serv   649
Priv-house-serv   149
Armed-Forces      9
Name: count, dtype: int64

```

Explore native_country variable

```
# check labels in native_country variable
```

```
df.native_country.unique()
```

```

↔ array(['United-States', 'Cuba', 'Jamaica', 'India', '?', 'Mexico',
       'South', 'Puerto-Rico', 'Honduras', 'England', 'Canada', 'Germany',
       'Iran', 'Philippines', 'Italy', 'Poland', 'Columbia', 'Cambodia',
       'Thailand', 'Ecuador', 'Laos', 'Taiwan', 'Haiti', 'Portugal',
       'Dominican-Republic', 'El-Salvador', 'France', 'Guatemala',
       'China', 'Japan', 'Yugoslavia', 'Peru',
       'Outlying-US(Guam-USVI-etc)', 'Scotland', 'Trinidad&Tobago',
       'Greece', 'Nicaragua', 'Vietnam', 'Hong', 'Ireland', 'Hungary',
       'Holand-Netherlands'], dtype=object)

```

```
# check frequency distribution of values in native_country variable
```

```
df.native_country.value_counts()
```

```

↔ native_country
United-States    29170
Mexico           643
?                583
Philippines      198
Germany          137
Canada           121
Puerto-Rico     114
El-Salvador      106
India            100
Cuba              95
England           90
Jamaica           81
South             80
China             75
Italy             73
Dominican-Republic 70
Vietnam           67
Guatemala        64
Japan             62
Poland           60

```

| | |
|----------------------------|----|
| Columbia | 59 |
| Taiwan | 51 |
| Haiti | 44 |
| Iran | 43 |
| Portugal | 37 |
| Nicaragua | 34 |
| Peru | 31 |
| France | 29 |
| Greece | 29 |
| Ecuador | 28 |
| Ireland | 24 |
| Hong | 20 |
| Cambodia | 19 |
| Trinidad&Tobago | 19 |
| Laos | 18 |
| Thailand | 18 |
| Yugoslavia | 16 |
| Outlying-US(Guam-USVI-etc) | 14 |
| Honduras | 13 |
| Hungary | 13 |
| Scotland | 12 |
| Holand-Netherlands | 1 |

Name: count, dtype: int64

We can see that there are 583 values encoded as ? in native_country variable. I will replace these ? with NaN.

```
# replace '?' values in native_country variable with `NaN`
df['native_country'].replace('?', np.NaN, inplace=True)
```

Check missing values in categorical variables again

```
df[categorical].isnull().sum()
```

```
workclass      1836
education       0
marital_status  0
occupation     1843
relationship    0
race            0
sex             0
native_country  583
income         0
dtype: int64
```

Now, we can see that workclass, occupation and native_country variable contains missing values.

Number of labels: cardinality The number of labels within a categorical variable is known as cardinality. A high number of labels within a variable is known as high cardinality. High cardinality may pose some serious problems in the machine learning model. So, I will check for high cardinality.

```
# check for cardinality in categorical variables
for var in categorical:
    print(var, ' contains ', len(df[var].unique()), ' labels')

workclass contains 9 labels
education contains 16 labels
marital_status contains 7 labels
occupation contains 15 labels
relationship contains 6 labels
race contains 5 labels
sex contains 2 labels
native_country contains 42 labels
income contains 2 labels
```

We can see that native_country column contains relatively large number of labels as compared to other columns. I will check for cardinality after train-test split

```
# find numerical variables
numerical = [var for var in df.columns if df[var].dtype!='O']

print('There are {} numerical variables\n'.format(len(numerical)))

print('The numerical variables are :', numerical)
```

There are 6 numerical variables

The numerical variables are : ['age', 'fnlwgt', 'education_num', 'capital_gain', 'capital_loss', 'hours_per_week']

view the numerical variables

```
df[numerical].head()
```

| | age | fnlwgt | education_num | capital_gain | capital_loss | hours_per_week |
|---|-----|--------|---------------|--------------|--------------|----------------|
| 0 | 39 | 77516 | 13 | 2174 | 0 | 40 |
| 1 | 50 | 83311 | 13 | 0 | 0 | 13 |
| 2 | 38 | 215646 | 9 | 0 | 0 | 40 |
| 3 | 53 | 234721 | 7 | 0 | 0 | 40 |
| 4 | 28 | 338409 | 13 | 0 | 0 | 40 |

Summary of numerical variables There are 6 numerical variables. These are given by age, fnlwgt, education_num, capital_gain, capital_loss and hours_per_week. All of the numerical variables are of discrete data type.

Missing values in numerical variables

check missing values in numerical variables

```
df[numerical].isnull().sum()
```

```
age      0
fnlwgt   0
education_num  0
capital_gain  0
capital_loss  0
hours_per_week  0
dtype: int64
```

We can see that all the 6 numerical variables do not contain missing values.

5. Declare feature vector and target variable

```
X = df.drop(['income'], axis=1)
```

```
y = df['income']
```

6. Split data into separate training and test set

split X and y into training and testing sets

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)
```

check the shape of X_train and X_test

```
X_train.shape, X_test.shape
```

```
((22792, 14), (9769, 14))
```

7. Feature Engineering

Feature Engineering is the process of transforming raw data into useful features that help us to understand our model better and increase its predictive power. I will carry out feature engineering on different types of variables.

First, I will display the categorical and numerical variables again separately.

check data types in X_train

```
X_train.dtypes
```



```

→ age          int64
  workclass    object
  fnlwgt       int64
  education    object
  education_num int64
  marital_status object
  occupation   object
  relationship object
  race         object
  sex          object
  capital_gain int64
  capital_loss int64
  hours_per_week int64
  native_country object
  dtype: object

```

```
# display categorical variables
```

```
categorical = [col for col in X_train.columns if X_train[col].dtypes == 'O']
```

```
categorical
```

```

→ ['workclass',
  'education',
  'marital_status',
  'occupation',
  'relationship',
  'race',
  'sex',
  'native_country']

```

```
# display numerical variables
```

```
numerical = [col for col in X_train.columns if X_train[col].dtypes != 'O']
```

```
numerical
```

```

→ ['age',
  'fnlwgt',
  'education_num',
  'capital_gain',
  'capital_loss',
  'hours_per_week']

```

Engineering missing values in categorical variables

```
# print percentage of missing values in the categorical variables in training set
```

```
X_train[categorical].isnull().mean()
```

```

→ workclass      0.055985
  education      0.000000
  marital_status  0.000000
  occupation      0.056072
  relationship    0.000000
  race           0.000000
  sex            0.000000
  native_country  0.018164
  dtype: float64

```

```
# print categorical variables with missing data
```

```

for col in categorical:
    if X_train[col].isnull().mean()>0:
        print(col, (X_train[col].isnull().mean()))

```

```

→ workclass 0.055984555984555984
  occupation 0.05607230607230607
  native_country 0.018164268164268166

```

```
# impute missing categorical variables with most frequent value
```

```

for df2 in [X_train, X_test]:
    df2['workclass'].fillna(X_train['workclass'].mode()[0], inplace=True)
    df2['occupation'].fillna(X_train['occupation'].mode()[0], inplace=True)
    df2['native_country'].fillna(X_train['native_country'].mode()[0], inplace=True)

```

```
# check missing values in categorical variables in X_train
```

```
X_train[categorical].isnull().sum()
```

```
workclass      0
education      0
marital_status 0
occupation     0
relationship    0
race           0
sex            0
native_country 0
dtype: int64
```

```
# check missing values in categorical variables in X_test
```

```
X_test[categorical].isnull().sum()
```

```
workclass      0
education      0
marital_status 0
occupation     0
relationship    0
race           0
sex            0
native_country 0
dtype: int64
```

I will check for missing values in X_train and X_test once again.

```
# check missing values in X_train
```

```
X_train.isnull().sum()
```

```
age           0
workclass     0
fnlwgt        0
education     0
education_num 0
marital_status 0
occupation    0
relationship  0
race          0
sex           0
capital_gain  0
capital_loss  0
hours_per_week 0
native_country 0
dtype: int64
```

```
# check missing values in X_test
```

```
X_test.isnull().sum()
```

```
age           0
workclass     0
fnlwgt        0
education     0
education_num 0
marital_status 0
occupation    0
relationship  0
race          0
sex           0
capital_gain  0
capital_loss  0
hours_per_week 0
native_country 0
dtype: int64
```

We can see that there are no missing values in X_train and X_test

Encode categorical variables

```
# print categorical variables
```

```
categorical
```

```
['workclass',
 'education',
 'marital_status',
 'occupation',
 'relationship',
 'race',
```

```
'sex',
'native_country']
```

```
X_train[categorical].head()
```

| | workclass | education | marital_status | occupation | relationship | race | sex | native_country |
|-------|-----------|--------------|--------------------|--------------|---------------|-------|--------|----------------|
| 32098 | Private | HS-grad | Married-civ-spouse | Craft-repair | Husband | White | Male | United-States |
| 25206 | State-gov | HS-grad | Divorced | Adm-clerical | Unmarried | White | Female | United-States |
| 23491 | Private | Some-college | Married-civ-spouse | Sales | Husband | White | Male | United-States |
| 12367 | Private | HS-grad | Never-married | Craft-repair | Not-in-family | White | Male | Guatemala |
| 7054 | Private | 7th-8th | Never-married | Craft-repair | Not-in-family | White | Male | Germany |

```
# import category encoders
!pip install category_encoders
import category_encoders as ce
```

```

Collecting category_encoders
  Downloading category_encoders-2.6.3-py2.py3-none-any.whl (81 kB)
    81.9/81.9 kB 3.1 MB/s eta 0:00:00
Requirement already satisfied: numpy>=1.14.0 in /usr/local/lib/python3.10/dist-packages (from category_encoders) (1.25.2)
Requirement already satisfied: scikit-learn>=0.20.0 in /usr/local/lib/python3.10/dist-packages (from category_encoders) (1.2.2)
Requirement already satisfied: scipy>=1.0.0 in /usr/local/lib/python3.10/dist-packages (from category_encoders) (1.11.4)
Requirement already satisfied: statsmodels>=0.9.0 in /usr/local/lib/python3.10/dist-packages (from category_encoders) (0.14.2)
Requirement already satisfied: pandas>=1.0.5 in /usr/local/lib/python3.10/dist-packages (from category_encoders) (2.0.3)
Requirement already satisfied: patsy>=0.5.1 in /usr/local/lib/python3.10/dist-packages (from category_encoders) (0.5.6)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.0.5->category_encoders) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.0.5->category_encoders) (2020.9)
Requirement already satisfied: tzdata>=2022.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.0.5->category_encoders) (2022.7)
Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packages (from patsy>=0.5.1->category_encoders) (1.16.0)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from scikit-learn>=0.20.0->category_encoders) (1.3.2)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn>=0.20.0->category_encoders) (3.2.0)
Requirement already satisfied: packaging>=21.3 in /usr/local/lib/python3.10/dist-packages (from statsmodels>=0.9.0->category_encoders) (23.1)
Installing collected packages: category_encoders
Successfully installed category_encoders-2.6.3

```

```
# encode remaining variables with one-hot encoding

encoder = ce.OneHotEncoder(cols=['workclass', 'education', 'marital_status', 'occupation', 'relationship',
                                'race', 'sex', 'native_country'])

X_train = encoder.fit_transform(X_train)

X_test = encoder.transform(X_test)
```

```
X_train.head()
```

| | age | workclass_1 | workclass_2 | workclass_3 | workclass_4 | workclass_5 | workclass_6 | workclass_7 | workclass_8 | fnlwgt | ... | nat: |
|--------------|-----|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------|-----|------|
| 32098 | 45 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 170871 | ... | |
| 25206 | 47 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 108890 | ... | |
| 23491 | 48 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 187505 | ... | |
| 12367 | 29 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 145592 | ... | |
| 7054 | 23 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 203003 | ... | |

5 rows × 105 columns


X_train.shape

(22792, 105)

We can see that from the initial 14 columns, we now have 113 columns.

Similarly, I will take a look at the `X_test` set.

```
X_test.head()
```



| | age | workclass_1 | workclass_2 | workclass_3 | workclass_4 | workclass_5 | workclass_6 | workclass_7 | workclass_8 | fnlwgt | ... | nat: |
|--------------|-----|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------|-----|------|
| 22278 | 27 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 177119 | ... | |
| 8950 | 27 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 216481 | ... | |
| 7838 | 25 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 256263 | ... | |
| 16505 | 46 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 147640 | ... | |
| 19140 | 45 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 172822 | ... | |

5 rows × 105 columns

X_test.shape

 (9769, 105)

We now have training and testing set ready for model building. Before that, we should map all the feature variables onto the same scale. It is called feature scaling. I will do it as follows.

8. Feature Scaling

```
cols = X_train.columns
```

```
from sklearn.preprocessing import RobustScaler
```

```
scaler = RobustScaler()
```


```
X_train = scaler.fit_transform(X_train)
```

```
X_test = scaler.transform(X_test)
```

```
X_train = pd.DataFrame(X_train, columns=[cols])
```

```
X_test = pd.DataFrame(X_test, columns=[cols])
```

```
X_train.head()
```



| | age | workclass_1 | workclass_2 | workclass_3 | workclass_4 | workclass_5 | workclass_6 | workclass_7 | workclass_8 | fnlwgt | ... | nati: |
|----------|-------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-----------|-----|-------|
| 0 | 0.40 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | -0.058906 | ... | |
| 1 | 0.50 | -1.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | -0.578076 | ... | |
| 2 | 0.55 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.080425 | ... | |
| 3 | -0.40 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | -0.270650 | ... | |
| 4 | -0.70 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.210240 | ... | |

5 rows × 105 columns


We now have X_train dataset ready to be fed into the Gaussian Naive Bayes classifier. I will do it as follows.

9. Model training

```
# train a Gaussian Naive Bayes classifier on the training set
from sklearn.naive_bayes import GaussianNB
```

```
# instantiate the model
gnb = GaussianNB()
```

```
# fit the model
gnb.fit(X_train, y_train)
```



▾ GaussianNB
 GaussianNB()

✓ 10. Predict the results

```
y_pred = gnb.predict(X_test)

y_pred
```

```
array(['<=50K', '<=50K', '>50K', ..., '>50K', '<=50K', '<=50K'],
      dtype='<U5')
```

✓ 11. Check accuracy score

```
from sklearn.metrics import accuracy_score

print('Model accuracy score: {0:0.4f}'.format(accuracy_score(y_test, y_pred)))
```

```
Model accuracy score: 0.8083
```

Here, `y_test` are the true class labels and `y_pred` are the predicted class labels in the test-set.

Compare the train-set and test-set accuracy Now, I will compare the train-set and test-set accuracy to check for overfitting.

```
y_pred_train = gnb.predict(X_train)

y_pred_train
```

```
array(['>50K', '<=50K', '>50K', ..., '<=50K', '>50K', '<=50K'],
      dtype='<U5')
```

```
print('Training-set accuracy score: {0:0.4f}'.format(accuracy_score(y_train, y_pred_train)))
```

```
Training-set accuracy score: 0.8067
```

Check for overfitting and underfitting

```
# print the scores on training and test set

print('Training set score: {:.4f}'.format(gnb.score(X_train, y_train)))

print('Test set score: {:.4f}'.format(gnb.score(X_test, y_test)))
```

```
Training set score: 0.8067
Test set score: 0.8083
```

The training-set accuracy score is 0.8067 while the test-set accuracy to be 0.8083. These two values are quite comparable. So, there is no sign of overfitting.

Compare model accuracy with null accuracy So, the model accuracy is 0.8083. But, we cannot say that our model is very good based on the above accuracy. We must compare it with the null accuracy. Null accuracy is the accuracy that could be achieved by always predicting the most frequent class.

So, we should first check the class distribution in the test set.

```
# check class distribution in test set

y_test.value_counts()
```

```
income
<=50K    7407
>50K     2362
Name: count, dtype: int64
```

We can see that the occurrences of most frequent class is 7407. So, we can calculate null accuracy by dividing 7407 by total number of occurrences.

```
# check null accuracy score

null_accuracy = (7407/(7407+2362))

print('Null accuracy score: {0:0.4f}'.format(null_accuracy))
```

Null accuracy score: 0.7582

We can see that our model accuracy score is 0.8083 but null accuracy score is 0.7582. So, we can conclude that our Gaussian Naive Bayes Classification model is doing a very good job in predicting the class labels.

Now, based on the above analysis we can conclude that our classification model accuracy is very good. Our model is doing a very good job in terms of predicting the class labels.

But, it does not give the underlying distribution of values. Also, it does not tell anything about the type of errors our classifier is making.

We have another tool called Confusion matrix that comes to our rescue.

12. Confusion matrix

A confusion matrix is a tool for summarizing the performance of a classification algorithm. A confusion matrix will give us a clear picture of classification model performance and the types of errors produced by the model. It gives us a summary of correct and incorrect predictions broken down by each category. The summary is represented in a tabular form.

Four types of outcomes are possible while evaluating a classification model performance. These four outcomes are described below:-

True Positives (TP) – True Positives occur when we predict an observation belongs to a certain class and the observation actually belongs to that class.

True Negatives (TN) – True Negatives occur when we predict an observation does not belong to a certain class and the observation actually does not belong to that class.

False Positives (FP) – False Positives occur when we predict an observation belongs to a certain class but the observation actually does not belong to that class. This type of error is called Type I error.

False Negatives (FN) – False Negatives occur when we predict an observation does not belong to a certain class but the observation actually belongs to that class. This is a very serious error and it is called Type II error.

```
# Print the Confusion Matrix and slice it into four pieces
```

```
from sklearn.metrics import confusion_matrix
```

```
cm = confusion_matrix(y_test, y_pred)
```

```
print('Confusion matrix\n\n', cm)
```

```
print('\nTrue Positives(TP) = ', cm[0,0])
```

```
print('\nTrue Negatives(TN) = ', cm[1,1])
```

```
print('\nFalse Positives(FP) = ', cm[0,1])
```

```
print('\nFalse Negatives(FN) = ', cm[1,0])
```

Confusion matrix

```
[[5999 1408]
 [ 465 1897]]
```

```
True Positives(TP) = 5999
```

```
True Negatives(TN) = 1897
```

```
False Positives(FP) = 1408
```

```
False Negatives(FN) = 465
```

The confusion matrix shows $5999 + 1897 = 7896$ correct predictions and $1408 + 465 = 1873$ incorrect predictions.

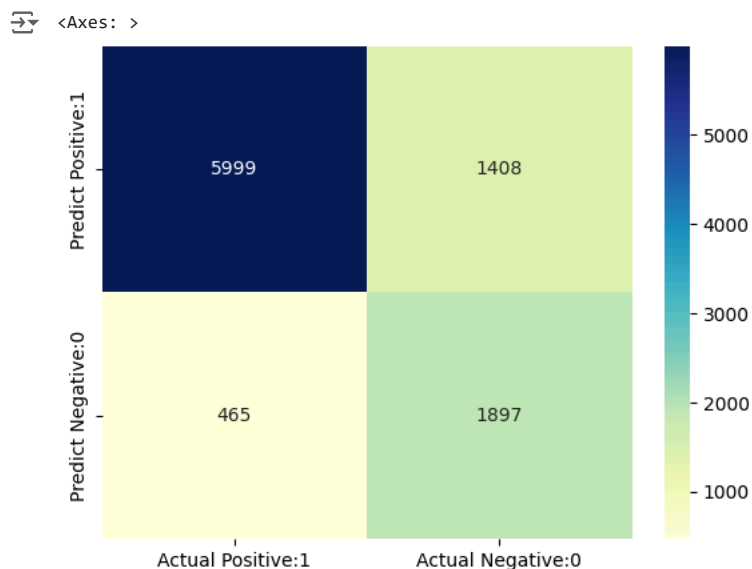
In this case, we have

True Positives (Actual Positive:1 and Predict Positive:1) - 5999 True Negatives (Actual Negative:0 and Predict Negative:0) - 1897 False Positives (Actual Negative:0 but Predict Positive:1) - 1408 (Type I error) False Negatives (Actual Positive:1 but Predict Negative:0) - 465 (Type II error)

```
# visualize confusion matrix with seaborn heatmap
```

```
cm_matrix = pd.DataFrame(data=cm, columns=['Actual Positive:1', 'Actual Negative:0'],
                        index=['Predict Positive:1', 'Predict Negative:0'])
```

```
sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='YlGnBu')
```



13. Classification metrics

Classification Report Classification report is another way to evaluate the classification model performance. It displays the precision, recall, f1 and support scores for the model. I have described these terms in later.

We can print a classification report as follows:-

```
from sklearn.metrics import classification_report
```

```
print(classification_report(y_test, y_pred))
```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| <=50K | 0.93 | 0.81 | 0.86 | 7407 |
| >50K | 0.57 | 0.80 | 0.67 | 2362 |
| accuracy | | | 0.81 | 9769 |
| macro avg | 0.75 | 0.81 | 0.77 | 9769 |
| weighted avg | 0.84 | 0.81 | 0.82 | 9769 |

Classification accuracy

```
TP = cm[0,0]
TN = cm[1,1]
FP = cm[0,1]
FN = cm[1,0]
```

```
# print classification accuracy
```

```
classification_accuracy = (TP + TN) / float(TP + TN + FP + FN)
```

```
print('Classification accuracy : {0:0.4f}'.format(classification_accuracy))
```

```
Classification accuracy : 0.8083
```

Classification error

```
# print classification error
```

```
classification_error = (FP + FN) / float(TP + TN + FP + FN)
```

```
print('Classification error : {0:0.4f}'.format(classification_error))
```

```
Classification error : 0.1917
```

Precision

Precision can be defined as the percentage of correctly predicted positive outcomes out of all the predicted positive outcomes. It can be given as the ratio of true positives (TP) to the sum of true and false positives (TP + FP).

So, Precision identifies the proportion of correctly predicted positive outcome. It is more concerned with the positive class than the negative class.

Mathematically, precision can be defined as the ratio of TP to (TP + FP)

```
# print precision score

precision = TP / float(TP + FP)

print('Precision : {0:0.4f}'.format(precision))
```

➞ Precision : 0.8099

Recall

Recall can be defined as the percentage of correctly predicted positive outcomes out of all the actual positive outcomes. It can be given as the ratio of true positives (TP) to the sum of true positives and false negatives (TP + FN). Recall is also called Sensitivity.

Recall identifies the proportion of correctly predicted actual positives.

Mathematically, recall can be given as the ratio of TP to (TP + FN).

```
recall = TP / float(TP + FN)

print('Recall or Sensitivity : {0:0.4f}'.format(recall))
```

➞ Recall or Sensitivity : 0.9281

True Positive Rate

True Positive Rate is synonymous with Recall.

```
true_positive_rate = TP / float(TP + FN)

print('True Positive Rate : {0:0.4f}'.format(true_positive_rate))
```

➞ True Positive Rate : 0.9281

False Positive Rate

```
false_positive_rate = FP / float(FP + TN)

print('False Positive Rate : {0:0.4f}'.format(false_positive_rate))
```

➞ False Positive Rate : 0.4260

Specificity

```
specificity = TN / (TN + FP)

print('Specificity : {0:0.4f}'.format(specificity))
```

➞ Specificity : 0.5740

f1-score

f1-score is the weighted harmonic mean of precision and recall. The best possible f1-score would be 1.0 and the worst would be 0.0. f1-score is the harmonic mean of precision and recall. So, f1-score is always lower than accuracy measures as they embed precision and recall into their computation. The weighted average of f1-score should be used to compare classifier models, not global accuracy.

✓ 14. Calculate class probabilities


```
# print the first 10 predicted probabilities of two classes- 0 and 1
```

```
y_pred_prob = gnb.predict_proba(X_test)[0:10]
```

```
y_pred_prob
```

```
array([[9.99999426e-01, 5.74152436e-07],
       [9.99687907e-01, 3.12093456e-04],
       [1.54405602e-01, 8.45594398e-01],
       [1.73624321e-04, 9.99826376e-01],
       [8.20121011e-09, 9.9999992e-01],
       [8.76844580e-01, 1.23155420e-01],
       [9.99999927e-01, 7.32876705e-08],
       [9.99993460e-01, 6.53998797e-06],
       [9.87738143e-01, 1.22618575e-02],
       [9.99999996e-01, 4.01886317e-09]])
```

Observations

In each row, the numbers sum to 1. There are 2 columns which correspond to 2 classes - <=50K and >50K.

Class 0 => <=50K - Class that a person makes less than equal to 50K.

Class 1 => >50K - Class that a person makes more than 50K.

Importance of predicted probabilities

We can rank the observations by probability of whether a person makes less than or equal to 50K or more than 50K. predict_proba process

Predicts the probabilities

Choose the class with the highest probability

Classification threshold level

There is a classification threshold level of 0.5.

Class 0 => <=50K - probability of salary less than or equal to 50K is predicted if probability < 0.5.

Class 1 => >50K - probability of salary more than 50K is predicted if probability > 0.5.

```
# store the probabilities in dataframe
```

```
y_pred_prob_df = pd.DataFrame(data=y_pred_prob, columns=['Prob of - <=50K', 'Prob of - >50K'])
```

```
y_pred_prob_df
```

| | Prob of - <=50K | Prob of - >50K | |
|---|-----------------|----------------|--|
| 0 | 9.999994e-01 | 5.741524e-07 | |
| 1 | 9.996879e-01 | 3.120935e-04 | |
| 2 | 1.544056e-01 | 8.455944e-01 | |
| 3 | 1.736243e-04 | 9.998264e-01 | |
| 4 | 8.201210e-09 | 1.000000e+00 | |
| 5 | 8.768446e-01 | 1.231554e-01 | |
| 6 | 9.999999e-01 | 7.328767e-08 | |
| 7 | 9.999935e-01 | 6.539988e-06 | |
| 8 | 9.877381e-01 | 1.226186e-02 | |
| 9 | 1.000000e+00 | 4.018863e-09 | |

Next steps: [Generate code with y_pred_prob_df](#) [View recommended plots](#)

```
# print the first 10 predicted probabilities for class 1 - Probability of >50K
```

```
gnb.predict_proba(X_test)[0:10, 1]
```

```
array([5.74152436e-07, 3.12093456e-04, 8.45594398e-01, 9.99826376e-01,
       9.9999992e-01, 1.23155420e-01, 7.32876705e-08, 6.53998797e-06,
       1.22618575e-02, 4.01886317e-09])
```

```
# store the predicted probabilities for class 1 - Probability of >50K
```

```
y_pred1 = gnb.predict_proba(X_test)[: , 1]
```

```
# adjust the font size
plt.rcParams['font.size'] = 12

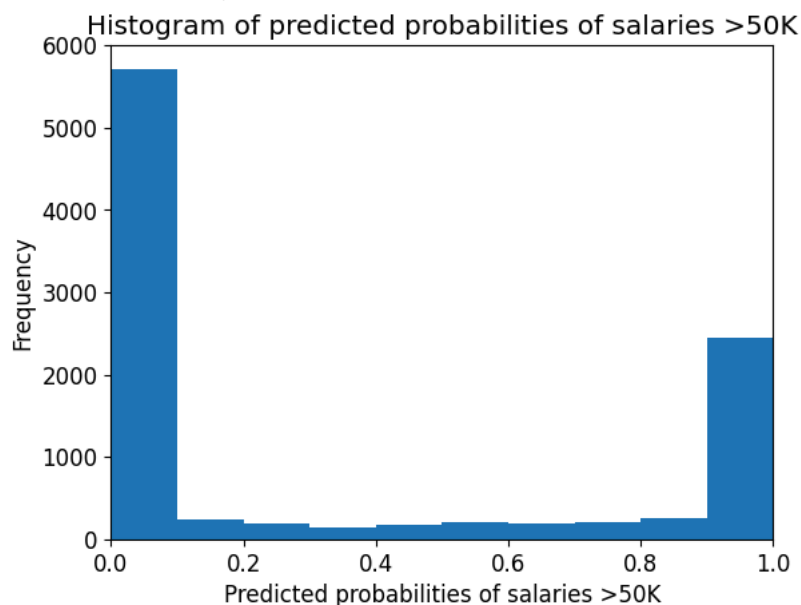
# plot histogram with 10 bins
plt.hist(y_pred1, bins = 10)

# set the title of predicted probabilities
plt.title('Histogram of predicted probabilities of salaries >50K')

# set the x-axis limit
plt.xlim(0,1)

# set the title
plt.xlabel('Predicted probabilities of salaries >50K')
plt.ylabel('Frequency')

Text(0, 0.5, 'Frequency')
```



Observations

We can see that the above histogram is highly positive skewed. The first column tell us that there are approximately 5700 observations with probability between 0.0 and 0.1 whose salary is $\leq 50K$. There are relatively small number of observations with probability > 0.5 . So, these small number of observations predict that the salaries will be $> 50K$. Majority of observations predict that the salaries will be $\leq 50K$.

✓ 14. ROC - AUC

ROC Curve Another tool to measure the classification model performance visually is ROC Curve. ROC Curve stands for Receiver Operating Characteristic Curve. An ROC Curve is a plot which shows the performance of a classification model at various classification threshold levels.

The ROC Curve plots the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold levels.

True Positive Rate (TPR) is also called Recall. It is defined as the ratio of TP to (TP + FN).

False Positive Rate (FPR) is defined as the ratio of FP to (FP + TN).

In the ROC Curve, we will focus on the TPR (True Positive Rate) and FPR (False Positive Rate) of a single point. This will give us the general performance of the ROC curve which consists of the TPR and FPR at various threshold levels. So, an ROC Curve plots TPR vs FPR at different classification threshold levels. If we lower the threshold levels, it may result in more items being classified as positive. It will increase both True Positives (TP) and False Positives (FP).

```
# plot ROC Curve

from sklearn.metrics import roc_curve

fpr, tpr, thresholds = roc_curve(y_test, y_pred1, pos_label = '>50K')

plt.figure(figsize=(6,4))

plt.plot(fpr, tpr, linewidth=2)

plt.plot([0,1], [0,1], 'k--' )

plt.rcParams['font.size'] = 12

plt.title('ROC curve for Gaussian Naive Bayes Classifier for Predicting Salaries')

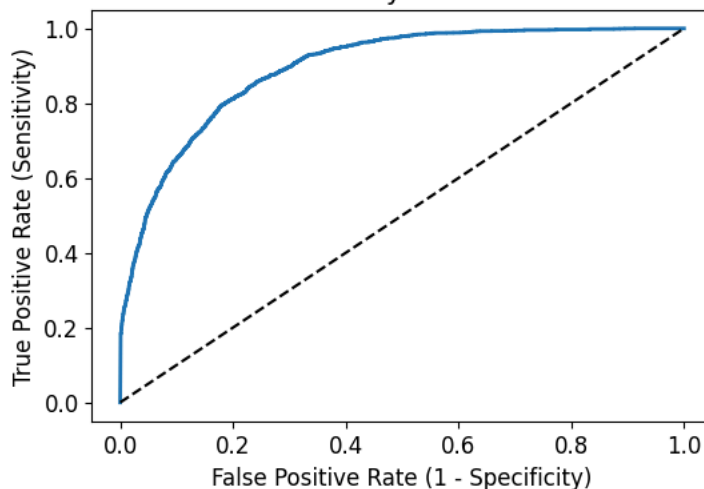
plt.xlabel('False Positive Rate (1 - Specificity)')

plt.ylabel('True Positive Rate (Sensitivity)')

plt.show()
```



ROC curve for Gaussian Naive Bayes Classifier for Predicting Salaries



ROC curve help us to choose a threshold level that balances sensitivity and specificity for a particular context.

ROC AUC

ROC AUC stands for Receiver Operating Characteristic - Area Under Curve. It is a technique to compare classifier performance. In this technique, we measure the area under the curve (AUC). A perfect classifier will have a ROC AUC equal to 1, whereas a purely random classifier will have a ROC AUC equal to 0.5.

So, ROC AUC is the percentage of the ROC plot that is underneath the curve

```
# compute ROC AUC

from sklearn.metrics import roc_auc_score

ROC_AUC = roc_auc_score(y_test, y_pred1)

print('ROC AUC : {:.4f}'.format(ROC_AUC))
```



ROC AUC : 0.8941

Interpretation¶ ROC AUC is a single number summary of classifier performance. The higher the value, the better the classifier.

ROC AUC of our model approaches towards 1. So, we can conclude that our classifier does a good job in predicting.

✓ 16. k-Fold Cross Validation

```
# Applying 10-Fold Cross Validation

from sklearn.model_selection import cross_val_score

scores = cross_val_score(gnb, X_train, y_train, cv = 10, scoring='accuracy')

We can summarize the cross-validation accuracy by calculating its mean

—
# compute Average cross-validation score

print('Average cross-validation score: {:.4f}'.format(scores.mean()))

↗ Average cross-validation score: 0.8063
```

Interpretation

Using the mean cross-validation, we can conclude that we expect the model to be around 80.63% accurate on average.

If we look at all the 10 scores produced by the 10-fold cross-validation, we can also conclude that there is a relatively small variance in the accuracy between folds, ranging from 81.35% accuracy to 79.64% accuracy. So, we can conclude that the model is independent of the particular folds used for training.

Our original model accuracy is 0.8083, but the mean cross-validation accuracy is 0.8063. So, the 10-fold cross-validation accuracy does not result in performance improvement for this model.

✓ 17. Results and conclusion

In this project, I build a Gaussian Naïve Bayes Classifier model to predict whether a person makes over 50K a year. The model yields a very good performance as indicated by the model accuracy which was found to be 0.8083. The training-set accuracy score is 0.8067 while the test-set accuracy to be 0.8083. These two values are quite comparable. So, there is no sign of overfitting. I have compared the model accuracy score which is 0.8083 with null accuracy score which is 0.7582. So, we can conclude that our Gaussian Naïve Bayes classifier model is doing a very good job in predicting the class labels. ROC AUC of our model approaches towards 1. So, we can conclude that our classifier does a very good job in predicting whether a person makes over 50K a year. Using the mean cross-validation, we can conclude that we expect the model to be around 80.63% accurate on average. If we look at all the 10 scores produced by the 10-fold cross-validation, we can also conclude that there is a relatively small variance in the accuracy between folds, ranging from 81.35% accuracy to 79.64% accuracy. So, we can conclude that the model is independent of the particular folds used for training. Our original model accuracy is 0.8083, but the mean cross-validation accuracy is 0.8063. So, the 10-fold cross-validation accuracy does not result in performance improvement for this model.

✓ 18. References

The work done in this project is inspired from following websites:-

<https://www.kaggle.com/datasets/qizarafzaal/adult-dataset>