



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

---

Experiment No. 6
Apply Boosting Algorithm on Adult Census Income Dataset and analyze the performance of the model
Date of Performance: 04-09-2023
Date of Submission: 09-10-2023



# Vidyavardhini's College of Engineering & Technology

## Department of Computer Engineering

---

**Aim:** Apply Boosting algorithm on Adult Census Income Dataset and analyze the performance of the model.

**Objective:** Apply Boosting algorithm on the given dataset and maximize the accuracy, Precision, Recall, F1 score.

### Theory:

Suppose that as a patient, you have certain symptoms. Instead of consulting one doctor, you choose to consult several. Suppose you assign weights to the value or worth of each doctor's diagnosis, based on the accuracies of previous diagnosis they have made. The final diagnosis is then a combination of the weighted diagnosis. This is the essence behind boosting.

Algorithm: Adaboost- A boosting algorithm—creates an ensemble of classifiers. Each one gives a weighted vote.

### Input:

- $D$ , a set of  $d$  class labelled training tuples
- $k$ , the number of rounds (one classifier is generated per round)
- a classification learning scheme

**Output:** A composite model

### Method

1. Initialize the weight of each tuple in  $D$  is  $1/d$
2. For  $i=1$  to  $k$  do // for each round
3. Sample  $D$  with replacement according to the tuple weights to obtain  $D_i$
4. Use training set  $D_i$  to derive a model  $M_i$
5. Compute  $\text{error}(M_i)$ , the error rate of  $M_i$
6.  $\text{Error}(M_i) = \sum w_j * \text{err}(X_j)$
7. If  $\text{Error}(M_i) > 0.5$  then
8. Go back to step 3 and try again
9. endif
10. for each tuple in  $D_i$  that was correctly classified do
11. Multiply the weight of the tuple by  $\text{error}(M_i)/(1-\text{error}(M_i))$
12. Normalize the weight of each tuple
13. end for



**To use the ensemble to classify tuple X**

1. Initialize the weight of each class to 0
2. for  $i=1$  to  $k$  do // for each classifier
3.  $w_i = \log((1 - \text{error}(M_i)) / \text{error}(M_i))$  // weight of the classifiers vote
4.  $C = M_i(X)$  // get class prediction for X from  $M_i$
5. Add  $w_i$  to weight for class C
6. end for
7. Return the class with the largest weight.

**Dataset:**

Predict whether income exceeds \$50K/yr based on census data. Also known as "Adult" dataset.

Attribute Information:

Listing of attributes:

>50K, <=50K.

age: continuous.

workclass: Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked.

fnlwgt: continuous.

education: Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 1st-4th, 10th, Doctorate, 5th-6th, Preschool.

education-num: continuous.

marital-status: Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse.



# Vidyavardhini's College of Engineering & Technology

## Department of Computer Engineering

---

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

relationship: Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried.

race: White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black.

sex: Female, Male.

capital-gain: continuous.

capital-loss: continuous.

hours-per-week: continuous.

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

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import io
from sklearn.metrics import accuracy_score, precision_score, f1_score, confusion_matrix, classification_report
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_squared_error
```

```
import os
for dirname, __, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

```
file = ('/content/adult.csv')
df = pd.read_csv(file)
```

```
print(df.head())
```

```

age workclass  fnlwgt      education  education.num marital.status \
0    90      ?    77053      HS-grad              9      Widowed
1    82  Private  132870      HS-grad              9      Widowed
2    66      ?   186061  Some-college             10      Widowed
3    54  Private  140359      7th-8th              4      Divorced
4    41  Private  264663  Some-college             10      Separated

      occupation  relationship    race    sex  capital.gain \
0              ?  Not-in-family  White  Female              0
1  Exec-managerial  Not-in-family  White  Female              0
2              ?    Unmarried  Black  Female              0
3  Machine-op-inspct    Unmarried  White  Female              0
4  Prof-specialty    Own-child   White  Female              0

capital.loss  hours.per.week  native.country  income
0          4356              40  United-States  <=50K
1          4356              18  United-States  <=50K
2          4356              40  United-States  <=50K
3          3900              40  United-States  <=50K
4          3900              40  United-States  <=50K
```

```
print(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
None
```

```

for i in df.columns:
    t = df[i].value_counts()
    index = list(t.index)
    print ("Count of ? in", i)
    for i in index:
        temp = 0
        if i == '?':
            print (t['?'])
            temp = 1
            break
```

```

if temp == 0:
    print ("0")

Count of ? in age
0
Count of ? in workclass
1836
Count of ? in fnlwgt
0
Count of ? in education
0
Count of ? in education.num
0
Count of ? in marital.status
0
Count of ? in occupation
1843
Count of ? in relationship
0
Count of ? in race
0
Count of ? in sex
0
Count of ? in capital.gain
0
Count of ? in capital.loss
0
Count of ? in hours.per.week
0
Count of ? in native.country
583
Count of ? in income
0

df=df.loc[(df['workclass'] != '?') & (df['native.country'] != '?')]
print(df.head())

```

	age	workclass	fnlwgt	education	education.num	marital.status	\
1	82	Private	132870	HS-grad	9	Widowed	
3	54	Private	140359	7th-8th	4	Divorced	
4	41	Private	264663	Some-college	10	Separated	
5	34	Private	216864	HS-grad	9	Divorced	
6	38	Private	150601	10th	6	Separated	

	occupation	relationship	race	sex	capital.gain	\
1	Exec-managerial	Not-in-family	White	Female	0	
3	Machine-op-inspct	Unmarried	White	Female	0	
4	Prof-specialty	Own-child	White	Female	0	
5	Other-service	Unmarried	White	Female	0	
6	Adm-clerical	Unmarried	White	Male	0	

	capital.loss	hours.per.week	native.country	income
1	4356	18	United-States	<=50K
3	3900	40	United-States	<=50K
4	3900	40	United-States	<=50K
5	3770	45	United-States	<=50K
6	3770	40	United-States	<=50K

```

df["income"] = [1 if i=='>50K' else 0 for i in df["income"]]
print(df.head())

```

	age	workclass	fnlwgt	education	education.num	marital.status	\
1	82	Private	132870	HS-grad	9	Widowed	
3	54	Private	140359	7th-8th	4	Divorced	
4	41	Private	264663	Some-college	10	Separated	
5	34	Private	216864	HS-grad	9	Divorced	
6	38	Private	150601	10th	6	Separated	

	occupation	relationship	race	sex	capital.gain	\
1	Exec-managerial	Not-in-family	White	Female	0	
3	Machine-op-inspct	Unmarried	White	Female	0	
4	Prof-specialty	Own-child	White	Female	0	
5	Other-service	Unmarried	White	Female	0	
6	Adm-clerical	Unmarried	White	Male	0	

	capital.loss	hours.per.week	native.country	income
1	4356	18	United-States	0
3	3900	40	United-States	0
4	3900	40	United-States	0
5	3770	45	United-States	0
6	3770	40	United-States	0

```
df_more=df.loc[df['income'] == 1]
print(df_more.head())
```

	age	workclass	fnlwtg	education	education.num	marital.status	\
7	74	State-gov	88638	Doctorate	16	Never-married	
10	45	Private	172274	Doctorate	16	Divorced	
11	38	Self-emp-not-inc	164526	Prof-school	15	Never-married	
12	52	Private	129177	Bachelors	13	Widowed	
13	32	Private	136204	Masters	14	Separated	

	occupation	relationship	race	sex	capital.gain	\
7	Prof-specialty	Other-relative	White	Female	0	
10	Prof-specialty	Unmarried	Black	Female	0	
11	Prof-specialty	Not-in-family	White	Male	0	
12	Other-service	Not-in-family	White	Female	0	
13	Exec-managerial	Not-in-family	White	Male	0	

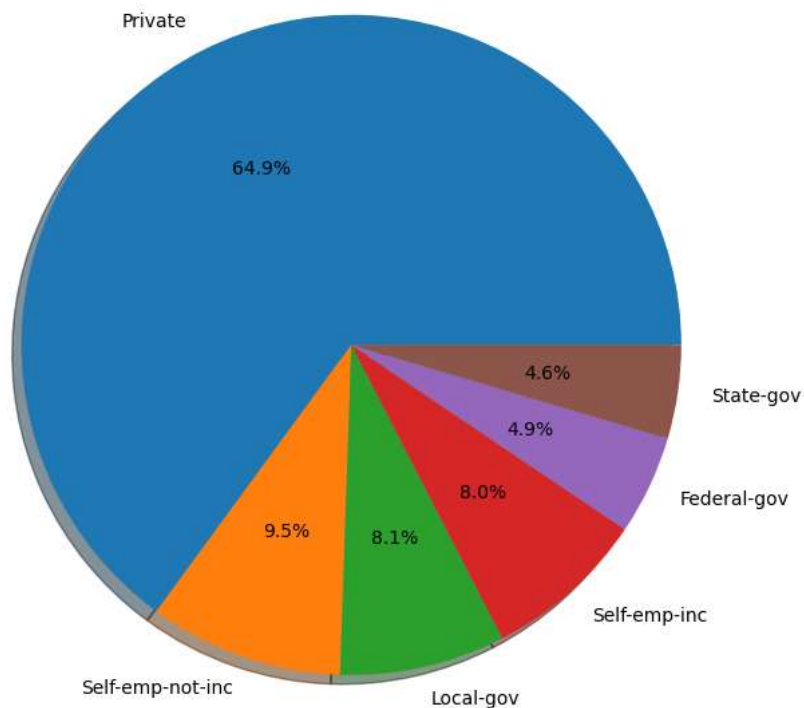
  

	capital.loss	hours.per.week	native.country	income
7	3683	20	United-States	1
10	3004	35	United-States	1
11	2824	45	United-States	1
12	2824	20	United-States	1
13	2824	55	United-States	1

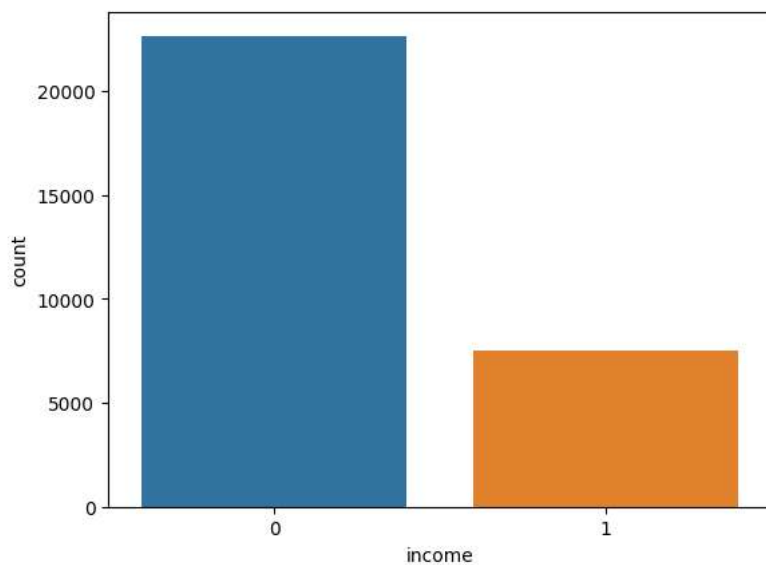
```
workclass_types = df_more['workclass'].value_counts()
labels = list(workclass_types.index)
aggregate = list(workclass_types)
print(workclass_types)
print(aggregate)
print(labels)
```

```
Private          4876
Self-emp-not-inc  714
Local-gov        609
Self-emp-inc      600
Federal-gov       365
State-gov         344
Name: workclass, dtype: int64
[4876, 714, 609, 600, 365, 344]
['Private', 'Self-emp-not-inc', 'Local-gov', 'Self-emp-inc', 'Federal-gov', 'State-gov']
```

```
plt.figure(figsize=(7,7))
plt.pie(aggregate, labels=labels, autopct='%1.1f%%', shadow = True)
plt.axis('equal')
plt.show()
```



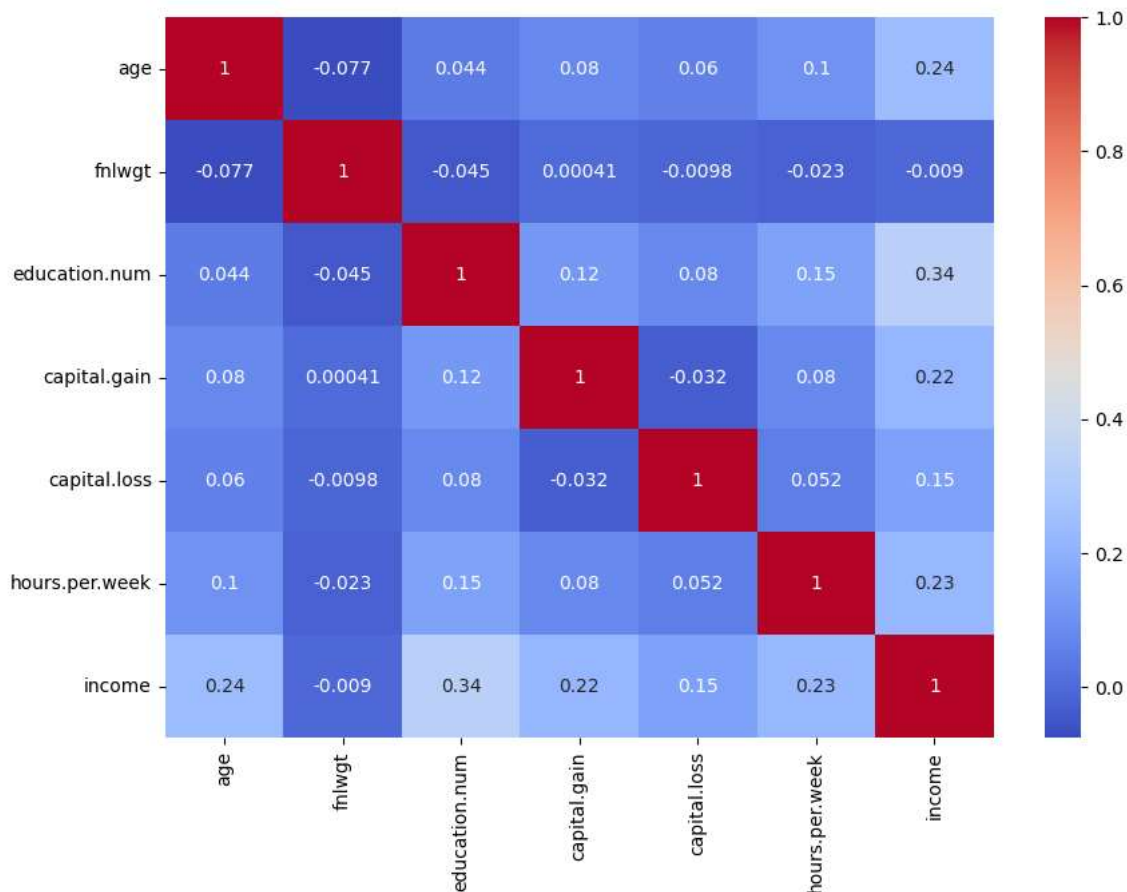
```
#Count plot on single categorical variable
sns.countplot(x='income', data = df)
plt.show()
df['income'].value_counts()
```



```
0    22661
1     7508
Name: income, dtype: int64
```

```
#Plot figsize
plt.figure(figsize=(10,7))
sns.heatmap(df.corr(), cmap='coolwarm', annot=True)
print(plt.show())
```

<ipython-input-13-6201d8194dba>:3: FutureWarning: The default value of numeric\_only in DataFrame.corr is deprecated. In a future version  
 sns.heatmap(df.corr(), cmap='coolwarm', annot=True)



None



```
plt.figure(figsize=(10,7))
sns.distplot(df['age'], color="red", bins=100)
plt.ylabel("Distribution", fontsize = 10)
plt.xlabel("Age", fontsize = 10)
plt.show()
```

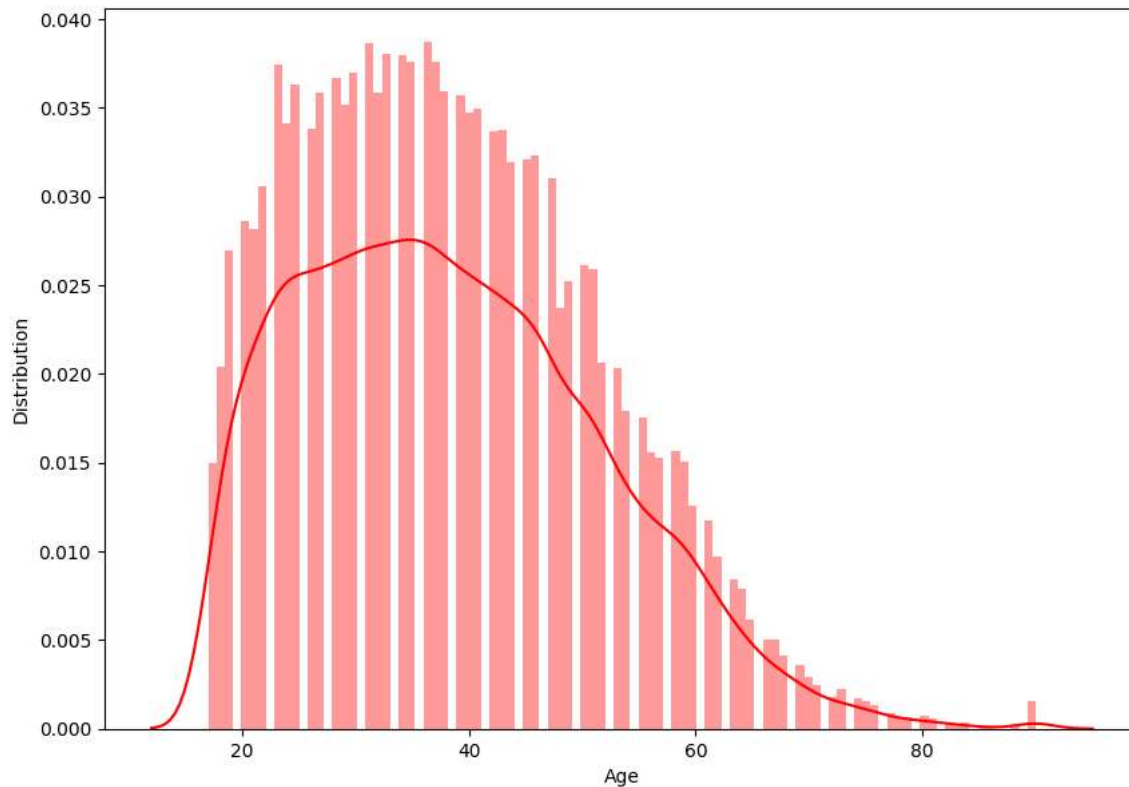
<ipython-input-14-1b72b8b67fa9>:2: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(df['age'], color="red", bins=100)
```



```
#To find distribution of categorical columns w.r.t income
fig, axes = plt.subplots(figsize=(20, 10))
```

```
plt.subplot(231)
sns.countplot(x='workclass',
              hue='income',
              data = df,
              palette="BuPu")
plt.xticks(rotation=90)
```

```
plt.subplot(232)
sns.countplot(x='marital.status',
              hue='income',
              data = df,
              palette="deep")
plt.xticks(rotation=90)
```

```
plt.subplot(233)
sns.countplot(x='education',
              hue='income',
              data = df,
              palette = "autumn")
plt.xticks(rotation=90)
```

```
plt.subplot(234)
sns.countplot(x='relationship',
              hue='income',
```

```

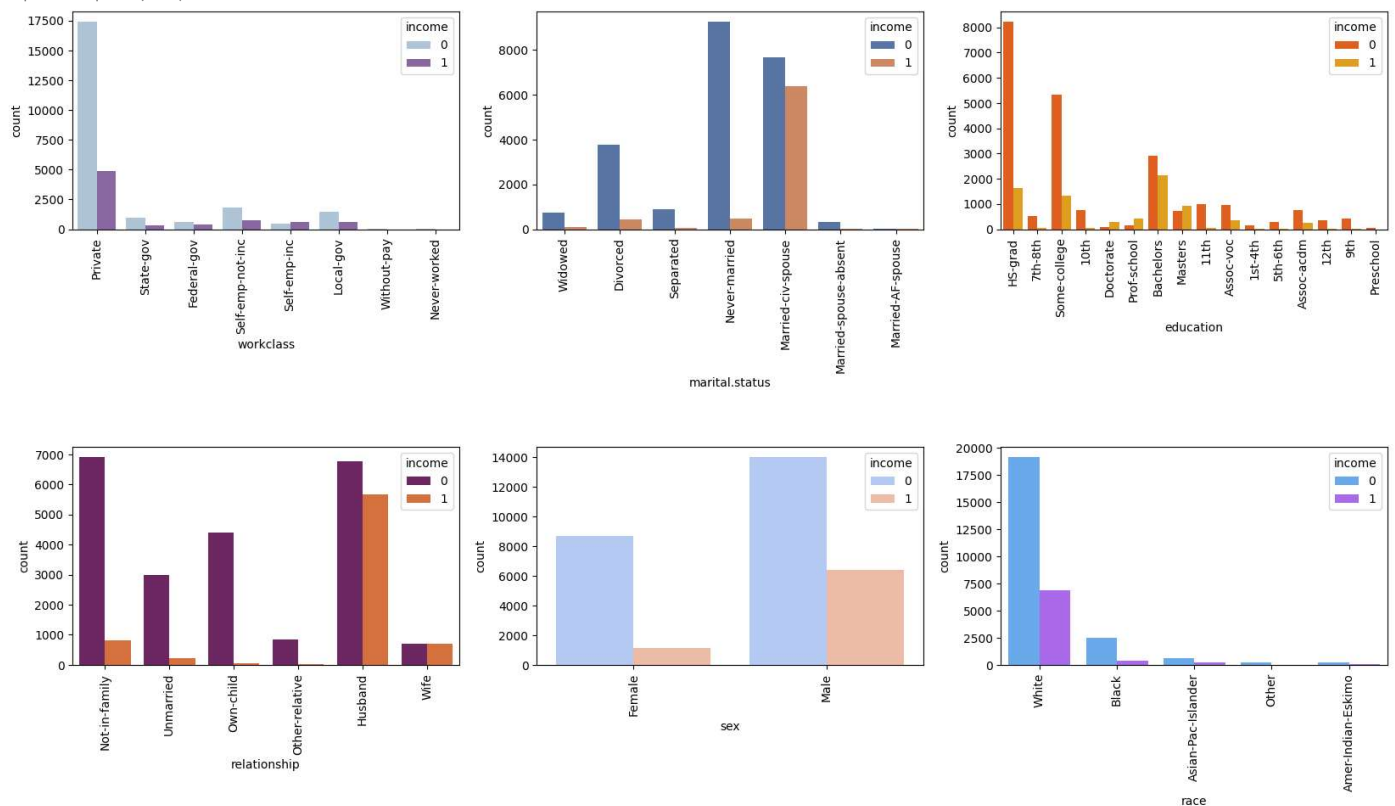
data = df,
palette = "inferno")
plt.xticks(rotation=90)

plt.subplot(235)
sns.countplot(x='sex',
              hue='income',
              data = df,
              palette = "coolwarm")
plt.xticks(rotation=90)

plt.subplot(236)
sns.countplot(x='race',
              hue='income',
              data = df,
              palette = "cool")
plt.xticks(rotation=90)
plt.subplots_adjust(hspace=1)
plt.show()

```

<ipython-input-15-f6a96c604872>:4: MatplotlibDeprecationWarning: Auto-removal of overlapping axes is deprecated since 3.6 and will be removed in a future version.



```

df1 = df.copy()

categorical_features = list(df1.select_dtypes(include=['object']).columns)
print(categorical_features)
df1

```

['workclass', 'education', 'marital.status', 'occupation', 'relationship', 'race', 'sex', 'native.country']

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.l
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	Not-in-family	White	Female	0	4
3	54	Private	140359	7th-8th	4	Divorced	Machine-op-inspct	Unmarried	White	Female	0	3
4	41	Private	264663	Some-college	10	Separated	Prof-specialty	Own-child	White	Female	0	3
5	34	Private	216864	HS-grad	9	Divorced	Other-service	Unmarried	White	Female	0	3
6	38	Private	150601	10th	6	Separated	Adm-clerical	Unmarried	White	Male	0	3
...	...	...	...	...	...	...	...	...	...	...	...	...
32556	22	Private	310152	Some-college	10	Never-married	Protective-serv	Not-in-family	White	Male	0	
32557	27	Private	257302	Assoc-acdm	12	Married-civ-spouse	Tech-support	Wife	White	Female	0	
32558	40	Private	154374	HS-grad	9	Married-civ-spouse	Machine-op-inspct	Husband	White	Male	0	

```
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
for feat in categorical_features:
    df1[feat] = le.fit_transform(df1[feat].astype(str))
df1
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss
1	82	3	132870	11	9	6	4	1	4	0	0	4356
3	54	3	140359	5	4	0	7	4	4	0	0	3900
4	41	3	264663	15	10	5	10	3	4	0	0	3900
5	34	3	216864	11	9	0	8	4	4	0	0	3770
6	38	3	150601	0	6	5	1	4	4	1	0	3770
...	...	...	...	...	...	...	...	...	...	...	...	...
32556	22	3	310152	15	10	4	11	1	4	1	0	0
32557	27	3	257302	7	12	2	13	5	4	0	0	0
32558	40	3	154374	11	9	2	7	0	4	1	0	0
32559	58	3	151910	11	9	6	1	4	4	0	0	0
32560	22	3	201490	11	9	4	1	3	4	1	0	0

30169 rows x 15 columns

```
X = df1.drop(columns = ['income'])
y = df1['income'].values

# Splitting the data set into train and test set
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)

print ("Train set size: ", X_train.shape)
print ("Test set size: ", X_test.shape)

Train set size: (21118, 14)
Test set size: (9051, 14)

from sklearn.ensemble import AdaBoostClassifier

# Train Adaboost Classifier
abc = AdaBoostClassifier(n_estimators = 300, learning_rate=1)
abc_model = abc.fit(X_train, y_train)

#Prediction
```

```

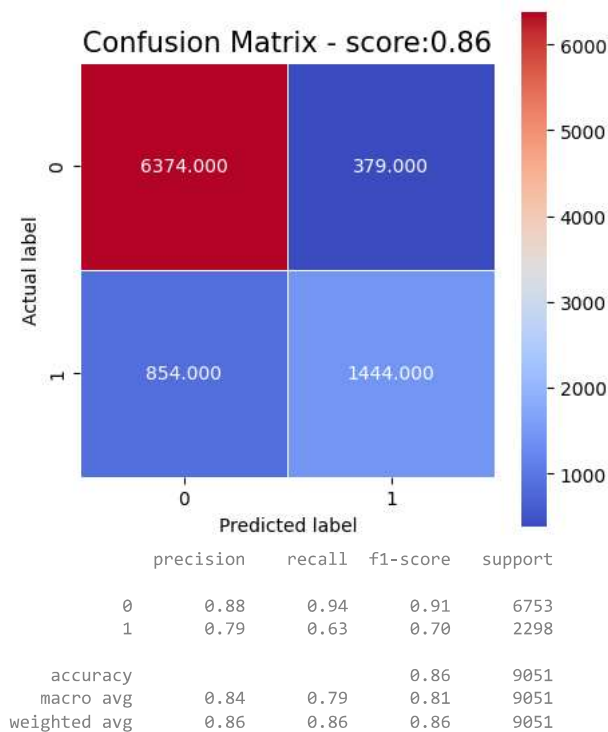
y_pred_abc = abc_model.predict(X_test)

print("Accuracy: ", accuracy_score(y_test, y_pred_abc))
print("F1 score :",f1_score(y_test, y_pred_abc, average='binary'))
print("Precision : ", precision_score(y_test, y_pred_abc))

Accuracy:  0.8637719588995691
F1 score : 0.7008007765105557
Precision : 0.7921009325287987

cm = confusion_matrix(y_test, y_pred_abc)
plt.figure(figsize=(5,5))
sns.heatmap(cm, annot=True, fmt=".3f", linewidths=.5, square = True, cmap = "coolwarm");
plt.ylabel('Actual label');
plt.xlabel('Predicted label');
plt.title('Confusion Matrix - score:' + str(round(accuracy_score(y_test, y_pred_abc), 2)), size = 15);
plt.show()
print(classification_report(y_test, y_pred_abc))

```



```

from sklearn.ensemble import GradientBoostingClassifier

#Training the model with gradient boosting
gbc = GradientBoostingClassifier(
    learning_rate = 0.1,
    n_estimators = 500,
    max_depth = 5,
    subsample = 0.9,
    min_samples_split = 100,
    max_features='sqrt',
    random_state=10)
gbc.fit(X_train,y_train)

# Predictions
y_pred_gbc = gbc.predict(X_test)

print("Accuracy : ",accuracy_score(y_test, y_pred_gbc))
print("F1 score : ", f1_score(y_test, y_pred_gbc, average = 'binary'))
print("Precision : ", precision_score(y_test, y_pred_gbc))

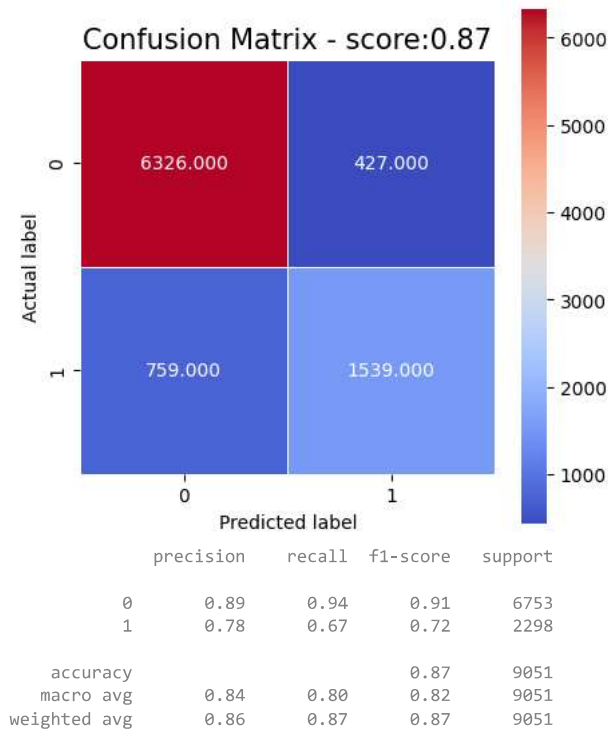
Accuracy : 0.8689647552756602
F1 score : 0.7218574108818011
Precision : 0.7828077314343845

```

```
rms = np.sqrt(mean_squared_error(y_test, y_pred_gbc))
print("RMSE for gradient boost: ", rms)
```

RMSE for gradient boost: 0.3619879068758235

```
cm = confusion_matrix(y_test, y_pred_gbc)
plt.figure(figsize=(5,5))
sns.heatmap(cm, annot = True, fmt=".3f", linewidths = 0.5, square = True, cmap = "coolwarm");
plt.ylabel('Actual label');
plt.xlabel('Predicted label');
plt.title('Confusion Matrix - score:' + str(round(accuracy_score(y_test, y_pred_gbc),2)), size = 15);
plt.show()
print(classification_report(y_test, y_pred_gbc))
```



```
import xgboost as xgb
from xgboost import XGBClassifier
```

```
#Training the model with gradient boosting
xgboost = XGBClassifier(learning_rate=0.01,
                        colsample_bytree = 0.4,
                        n_estimators=1000,
                        max_depth=20,
                        gamma=1)
```

```
xgboost_model = xgboost.fit(X_train, y_train)
```

```
# Predictions
y_pred_xgboost = xgboost_model.predict(X_test)
```

```
print("Accuracy : ",accuracy_score(y_test, y_pred_xgboost))
print("F1 score : ", f1_score(y_test, y_pred_xgboost, average = 'binary'))
print("Precision : ", precision_score(y_test, y_pred_xgboost))
```

Accuracy : 0.8655397193680257  
 F1 score : 0.7090604829070045  
 Precision : 0.786737400530504

```
rms = np.sqrt(mean_squared_error(y_test, y_pred_xgboost))
print("RMSE for xgboost: ", rms)
```

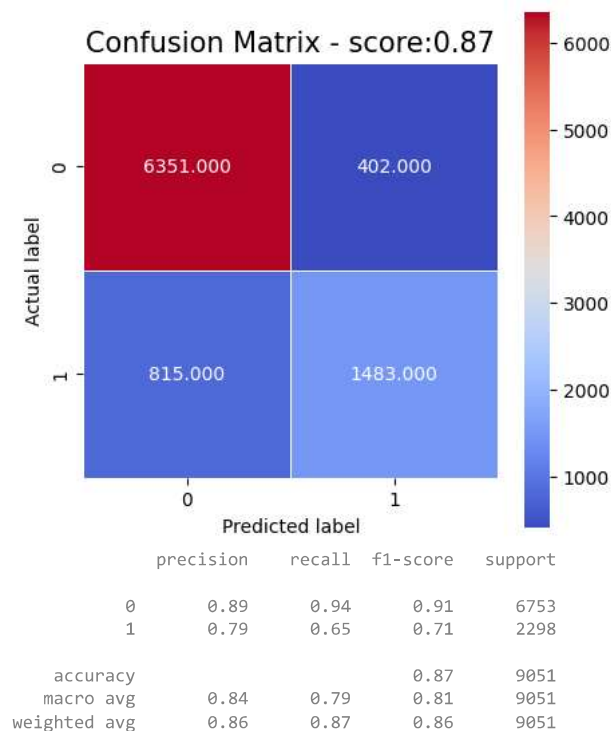
RMSE for xgboost: 0.3666882608319693

```
cm = confusion_matrix(y_test, y_pred_xgboost)
plt.figure(figsize=(5,5))
```

```

sns.heatmap(cm, annot=True, fmt=".3f", linewidths=.5, square = True, cmap = "coolwarm");
plt.ylabel('Actual label');
plt.xlabel('Predicted label');
plt.title('Confusion Matrix - score:'+str(round(accuracy_score(y_test, y_pred_xgboost),2)), size = 15);
plt.show()
print(classification_report(y_test,y_pred_xgboost))

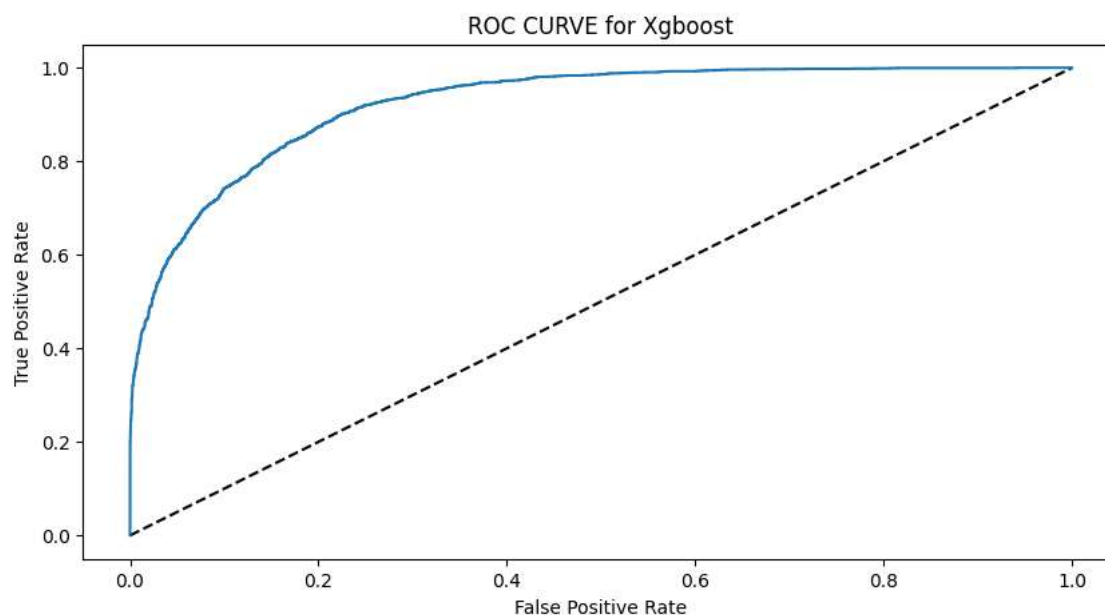
```



```

from sklearn.metrics import roc_curve
fpr, tpr, thresholds = roc_curve(y_test, xgboost.predict_proba(X_test)[:,:1])
plt.figure(figsize = (10,5))
plt.plot([0,1],[0,1], 'k--')
plt.plot(fpr, tpr)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC CURVE for Xgboost')
plt.show()

```







# Vidyavardhini's College of Engineering & Technology

## Department of Computer Engineering

---

### Conclusion:

1. The GradientBoostingClassifier has the highest accuracy (0.8690) and F1 score (0.7219) among the three classifiers, indicating that it performs the best on this dataset.
2. The AdaBoostClassifier and XGBClassifier also have good performance but slightly lower than the GradientBoostingClassifier in terms of accuracy and F1 score.
3. All three classifiers have relatively high precision, indicating that when they predict a positive class (1), they are often correct.
4. The recall values vary among the classifiers, with the GradientBoostingClassifier having the highest recall for the positive class (1), indicating that it correctly identifies more positive cases.
5. The F1 score, which balances precision and recall, shows how well the classifiers perform in classifying both classes.

### Comparison between boosting algorithms and random forest classifier:

1. The boosting algorithms (AdaBoost, Gradient Boosting, and XGBoost) generally outperform the Random Forest Classifier in terms of accuracy, precision, and F1 score.
2. The Random Forest Classifier performs reasonably well, with an accuracy of around 85% and a balanced F1 score. However, it falls slightly behind the boosting algorithms.
3. The boosting algorithms tend to have higher precision and recall for the positive class (income > 50K), indicating better ability to correctly classify high-income individuals.

All these models perform reasonably well on the Adult Census Income Dataset; the boosting algorithms, especially Gradient Boosting, appear to provide slightly better results in terms of accuracy and F1 score compared to the Random Forest Classifier.