

```

import pandas as pd
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt
from sklearn.preprocessing import LabelEncoder
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split, cross_val_score, KFold, GridSearchCV
from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
# import scikitplot as skplt

```

```
dataset=pd.read_csv("/content/adult.csv")
```

```

print(dataset.isnull().sum())
print(dataset.dtypes)

```

```

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
income       0
dtype: int64
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
income       object
dtype: object

```

```
dataset.head()
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	race	sex	capital.gain	capital.loss
0	90	?	77053	HS-grad	9	Widowed	?	Not-in-family	White	Female	0	4356
1	82	Private	132870	HS-grad	9	Widowed	Exec-managerial	Not-in-family	White	Female	0	4356
2	66	?	186061	Some-college	10	Widowed	?	Unmarried	Black	Female	0	4356
3	54	Private	140359	7th-8th	4	Divorced	Machine-op-inspct	Unmarried	White	Female	0	3900
4	41	Private	264663	Some-college	10	Separated	Prof-specialty	Own-child	White	Female	0	3900

```

#removing '?' containing rows
dataset = dataset[(dataset != '?').all(axis=1)]
#label the income objects as 0 and 1
dataset['income']=dataset['income'].map({'<=50K': 0, '>50K': 1})

```

```

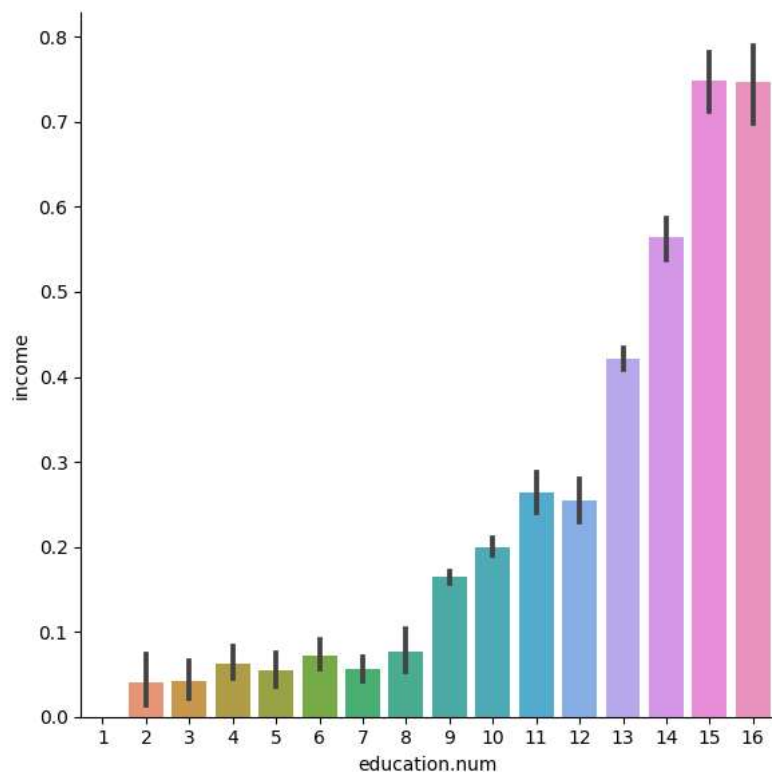
<ipython-input-6-39ed73805135>:4: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.

```

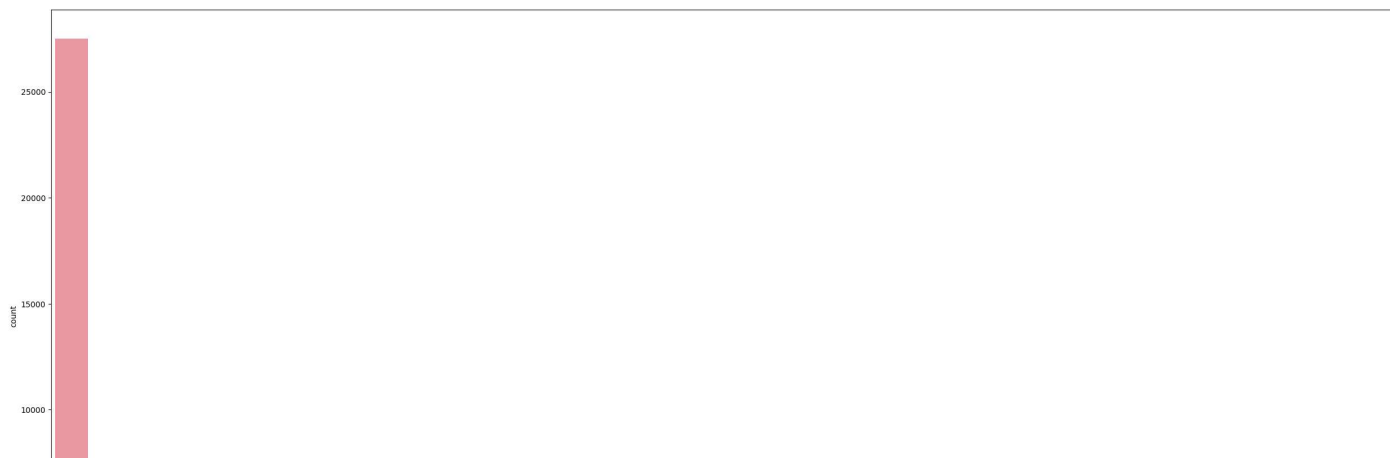
Try using `.loc[row_indexer,col_indexer] = value` instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
`dataset['income']=dataset['income'].map({'<=50K': 0, '>50K': 1})`

```
sns.catplot(x='education.num',y='income',data=dataset,kind='bar',height=6)
plt.show()
```



```
#explore which country do most people belong
plt.figure(figsize=(38,14))
sns.countplot(x='native.country',data=dataset)
plt.show()
```



```
for column in dataset:
    enc=LabelEncoder()
    if dataset.dtypes[column]==np.object:
        dataset[column]=enc.fit_transform(dataset[column])
```

<ipython-input-10-5d7d7fe4d7c0>:3: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warn_

Deprecated in NumPy 1.20; for more details and guidance: <https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations>

<ipython-input-10-5d7d7fe4d7c0>:3: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warn_

Deprecated in NumPy 1.20; for more details and guidance: <https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations>

<ipython-input-10-5d7d7fe4d7c0>:3: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warn_

```

if dataset.dtypes[column]==np.object:
<ipython-input-10-5d7d7fe4d7c0>:3: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warn
Depreciated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
    if dataset.dtypes[column]==np.object:
<ipython-input-10-5d7d7fe4d7c0>:3: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warn
Depreciated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
    if dataset.dtypes[column]==np.object:
<ipython-input-10-5d7d7fe4d7c0>:3: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warn
Depreciated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
    if dataset.dtypes[column]==np.object:
<ipython-input-10-5d7d7fe4d7c0>:3: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warn
Depreciated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
    if dataset.dtypes[column]==np.object:
<ipython-input-10-5d7d7fe4d7c0>:3: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warn
Depreciated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
    if dataset.dtypes[column]==np.object:
<ipython-input-10-5d7d7fe4d7c0>:3: DeprecationWarning: `np.object` is a deprecated alias for the builtin `object`. To silence this warn
Depreciated in NumPy 1.20; for more details and guidance: https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations
    if dataset.dtypes[column]==np.object:

```

age	1.00	0.08	-0.08	-0.00	0.04	-0.28	-0.01	-0.25	0.02	0.08	0.08	0.06	0.10	-0.00	0.24
workclass	0.08	1.00	-0.03	0.02	0.04	-0.03	0.02	-0.07	0.04	0.07	0.04	0.01	0.05	0.01	0.02
fnlwtg	-0.08	-0.03	1.00	-0.03	-0.04	0.03	0.00	0.01	-0.02	0.03	0.00	-0.01	-0.02	-0.07	-0.01
education	-0.00	0.02	-0.03	1.00	0.35	-0.04	-0.04	-0.01	0.01	-0.03	0.03	0.02	0.06	0.08	0.08
education.num	0.04	0.04	-0.04	0.35	1.00	-0.06	0.09	-0.09	0.03	0.01	0.12	0.08	0.15	0.09	0.34
marital.status	-0.28	-0.03	0.03	-0.04	-0.06	1.00	0.02	0.18	-0.07	-0.12	-0.04	-0.04	-0.19	-0.03	-0.19
occupation	-0.01	0.02	0.00	-0.04	0.09	0.02	1.00	-0.05	0.00	0.06	0.02	0.01	0.02	-0.00	0.05
relationship	-0.25	-0.07	0.01	-0.01	-0.09	0.18	-0.05	1.00	-0.12	-0.58	-0.06	-0.06	-0.26	-0.01	-0.25
race	0.02	0.04	-0.02	0.01	0.03	-0.07	0.00	-0.12	1.00	0.09	0.01	0.02	0.05	0.12	0.07
sex	0.08	0.07	0.03	-0.03	0.01	-0.12	0.06	-0.58	0.09	1.00	0.05	0.05	0.23	0.00	0.22
capital.gain	0.08	0.04	0.00	0.03	0.12	-0.04	0.02	-0.06	0.01	0.05	1.00	-0.03	0.08	0.01	0.22
capital.loss	0.06	0.01	-0.01	0.02	0.08	-0.04	0.01	-0.06	0.02	0.05	-0.03	1.00	0.05	0.01	0.15
hours.per.week	0.10	0.05	-0.02	0.06	0.15	-0.19	0.02	-0.26	0.05	0.23	0.08	0.05	1.00	0.01	0.23
native.country	-0.00	0.01	-0.07	0.08	0.09	-0.03	-0.00	-0.01	0.12	0.00	0.01	0.01	0.01	1.00	0.02

income 0.24 0.02 0.01 0.08 0.34 0.19 0.05 0.25 0.07 0.22 0.22 0.15 0.23 0.02 1.00

```
print(dataset.head())
```

	age	workclass	education.num	marital.status	race	sex	capital.gain	\
1	82	2	9	6	4	0	0	
3	54	2	4	0	4	0	0	
4	41	2	10	5	4	0	0	
5	34	2	9	0	4	0	0	
6	38	2	6	5	4	1	0	
	capital.loss	hours.per.week	income					
1	4356	18	0					
3	3900	40	0					

4	3900	40	0
5	3770	45	0
6	3770	40	0

```
X=dataset.iloc[:,0:-1]
y=dataset.iloc[:, -1]
print(X.head())
print(y.head())
x_train,x_test,y_train,y_test=train_test_split(X,y,test_size=0.33,shuffle=False)
```

	age	workclass	education.num	marital.status	race	sex	capital.gain \
1	82	2	9	6	4	0	0
3	54	2	4	0	4	0	0
4	41	2	10	5	4	0	0
5	34	2	9	0	4	0	0
6	38	2	6	5	4	1	0

	capital.loss	hours.per.week
1	4356	18
3	3900	40
4	3900	40
5	3770	45
6	3770	40

1	0
3	0
4	0
5	0
6	0

Name: income, dtype: int64

```
import xgboost as xgb
xgb.__version__
```

'2.0.0'

```
dmat=xgb.DMatrix(x_train,y_train)
test_dmat=xgb.DMatrix(x_test)
```

```
from skopt import BayesSearchCV
import warnings
warnings.filterwarnings('ignore', message='The objective has been evaluated at this point before.')
```

```
params={'min_child_weight': (0, 10),
        'max_depth': (0, 30),
        'subsample': (0.5, 1.0, 'uniform'),
        'colsample_bytree': (0.5, 1.0, 'uniform'),
        'n_estimators': (50,100),
        'reg_lambda': (1,100, 'log-uniform'),
        }
```

```
bayes=BayesSearchCV(estimator=xgb.XGBClassifier(objective='binary:logistic',eval_metric='error',eta=0.1),search_spaces=params,n_iter=50,scori
res=bayes.fit(x_train,y_train)
print(res.best_params_)
print(res.best_score_)
```

{'colsample_bytree': 1.0, 'max_depth': 30, 'min_child_weight': 0, 'n_estimators': 50, 'reg_lambda': 100.0, 'subsample': 0.5} 0.7750395882818686

```
final_clf=xgb.train(params=final_p,dtrain=dmat,num_boost_round=837)
pred=final_clf.predict(test_dmat)
print(pred)
pred[pred > 0.5 ] = 1
pred[pred <= 0.5] = 0
print(pred)
print(accuracy_score(y_test,pred)*100)
```

[8.5713279e-01 7.6487666e-01 5.4465812e-01 ... 2.9337847e-01 3.4949776e-02 5.5424345e-04] [1. 1. 1. ... 0. 0. 0.] 85.01105083383564

