

In []:

```

1  #####
2  #
3  #           Problem : To predict if a person makes >50K using XGBoost
4  #
5  #
6  #####
7  # Step 0: To Load the Libraries, packages being used
8
9  import numpy as np
10 import pandas as pd
11 from sklearn.preprocessing import LabelEncoder
12 import seaborn as sns
13 from sklearn.model_selection import train_test_split
14 import xgboost as xgb
15 from sklearn.metrics import roc_auc_score, confusion_matrix
16 from sklearn.model_selection import KFold, cross_val_score
17 import matplotlib.pyplot as plt
18 %matplotlib inline
19
20
21 # To Load the Data
22 train_set = pd.read_csv("https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data", header = None)
23
24 test_set = pd.read_csv("https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.test", skiprows = 1, header = None)
25
26 col_labels = ['age', 'workclass', 'fnlwgt', 'education', 'education_num', 'marital_status', 'occupation', 'relationship', 'sex', 'sex_num', 'race', 'race_num', 'ethnicity', 'ethnicity_num', 'income']
27 train_set.columns = col_labels
28 test_set.columns = col_labels
29
30 print(" The Training Data Set are as : \n", '-'*80)
31
32 print(train_set.head(2))
33 print('-'*80)
34
35 print(" The Test Data Set are as : \n", '-'*80)
36
37 print(test_set.head(2))
38
39 print( "\n The number of rows in training and test set \n", '-'*80)
40 print(train_set.shape)
41 print(test_set.shape)

```

```
42 train_rowcnt = train_set.shape[0]
43 test_rowcnt = test_set.shape[0]
44 train_percent = (train_rowcnt/(train_rowcnt+test_rowcnt))
45
46 print( " The Training percentage of rows is : ",round(train_percent,2)*100)
47
```

```
In [ ]: 1 # Step1 : Analysing the train Data by seeing info, describe, shape , nulls functions
2
3 print(" To combine the Train and Test DataSets to analyse together the Data ")
4 df = pd.concat([train_set,test_set], ignore_index= True)
5 print(df.shape)
6
7 print('-'*80)
8 print(df.head())
9
10 print(" The Total numbr of rows and columns in the DataSet \n",'-'*80)
11 print(df.shape)
12
13 print(" Analyse the Data by using info, describe and check for nulls, categorical fields ")
14 print('-'*80)
15
16 print(df.info())
17 print('-'*80)
18
19 print(df.describe())
20 print('-'*80)
21
22 print(df.isnull().sum())
23 print('-'*80)
```

```
In [ ]: 1 # Step2: Preprocessing the Data
2
3 # We observe , there are no nulls but lot of categorical fields which need to made numeric
4 #['workclass', 'education', 'marital_status', 'occupation', 'relationship', 'race', 'sex', 'native_country', 'wage_class
5
6 # using LabelEncoder to convert the Categorical fields to numeric in both train and test set
7 print(" The Categorical Fields and their unique values \n", '-'*80)
8 print(df.workclass.unique())
9 print(df.education.unique())
10 print(df.marital_status.unique())
11
12 print(df.occupation.unique())
13 print(df.relationship.unique())
14 print(df.race.unique())
15 print(df.sex.unique())
16 print(df.native_country.unique())
17 print(df.wage_class.unique())
18
19 # Observing the Data , we see there are spaces as well as '?' character in some fields , using function
20 # trimAllColumns to remove spaces
21
22 def trimAllColumns(df):
23     trimStrings = lambda x: x.strip() if type(x) is str else x
24     return df.applymap(trimStrings)
25
26 df = trimAllColumns(df)
27 df = df.replace('?', np.nan)
28
29 # also in the wage_Class we need only see if greater than 50K or not , so correctly some rows
30 df = df.replace('>50K.', '>50K')
31 df = df.replace('<=50K.', '<=50K')
32
33
```

```
In [ ]: 1 # since the number of nulls is very less as compared to size of the total rows , so we can ignore the nulls
2 # or drop it
3
4 df = df.dropna()
5
```

```
In [ ]: 1 # Changing the Categorical Values to numeric values
2
3 le = LabelEncoder()
4
5 cat_features = ['workclass', 'education', 'marital_status', 'occupation', 'relationship', 'race', 'sex', 'native_country', 'v
6
7
8 for att in cat_features:
9     df[att] = le.fit_transform(df[att].astype(str))
10
11 print(" The Categorical fields changed to numerical by Labelencoder \n", '-'*80)
12 print(df.head())
13
```

```
In [ ]: 1 # Step3 : Visualisation of train Data to see relationship with Wage_class attribute
2 # here instead of taking all the rows , a sample size of 10000 records is used
3
4
5 features = df.columns
6 y = 'wage_class'
7
8 print(" Visualising the Data to see the relationship with wage_class")
9 column_cnt = len(features)
10 i = 0
11 while i < column_cnt:
12     sns.pairplot(df.sample(10000), x_vars=features[i:i+4], y_vars=y, kind="reg")
13     i = i+4
14
```

```
In [ ]: 1 # The Full DataSets
        2
        3 X = df
        4 X = X.drop('wage_class',axis=1)
        5 y = df['wage_class']
        6 y = y.values.reshape(-1,1)
```

```
In [ ]: 1 #! pip install xgboost
        2 #Step4: Splitting the combine Dataset after preprocessing based on the previous training set size
        3
        4 X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=(1-train_perct), random_state=2)
        5
        6 # X_train = df[:train_set.shape[0]]
        7 # X_test = df[train_set.shape[0]:]
        8 print(" The number of rows amd columns in the training and test set now are :\n")
        9 print(X_train.shape)
       10 print(X_test.shape)
       11 print(y_train.shape)
       12
       13 print(" The training set : \n", '-'*80)
       14 print(X_train.head(2))
       15
       16 print(" The test set : \n", '-'*80)
       17 print(X_test.head(2))
       18
```

```
In [ ]: 1 # Step4: XgBoost model creation from sklearn library
2 from sklearn.metrics import accuracy_score
3 from xgboost.sklearn import XGBClassifier
4 from sklearn.model_selection import KFold, cross_val_score
5 import warnings
6
7
8 # to remove the warnings coming in the output
9 if __name__ == '__main__':
10     warnings.filterwarnings(action='ignore', category=DeprecationWarning)
11
12
13 # setting the parameters
14 params = {
15     'objective': 'binary:logistic',
16     'max_depth': 3,
17     'learning_rate': 0.1,
18     'silent': 1,
19     'n_estimators': 26,
20     'eval_metric': ['auc', 'logloss']
21 }
22
23
24 # Seeing the Performance For Xgboost Classifier with customised paramters
25 # Sklearn Xgboost Classifier 1
26
27 bst = XGBClassifier(**params)
28
29 y_train = y_train.ravel()
30
31 # training the Classifier with the training set
32 bst.fit(X_train, y_train)
33 #, eval_set=[(X_train, y_train), (X_test, y_test)]
34
35 print(" Xgboost Classifier 1 with Customised Parameters \n", '-'*80)
36
37 print(bst)
38
39 # predicting the labels for the test set
40 preds = bst.predict(X_test)
41
```

```
42 preds = pd.Series(preds)
43 print(preds.shape)
44 print(y_test.shape)
45
46 # y_test = y_test.values.reshape(-1,1)
47 preds = preds.values.reshape(-1,1)
48
49 correct=0
50 for i in range(len(preds)):
51     if (y_test[i] == preds[i]):
52         correct += 1
53 if y_test.size > 0 :
54     acc = accuracy_score(y_test, preds)
55
56 print('Predicted correctly: {0}/{1}'.format(correct, len(preds)))
57 print('Accuracy calculated manually: ', correct*100/len(preds), '%')
58 print('Accuracy calculated by accuracy score function: ', acc*100, '%')
59 print('Error: {0:.4f}'.format((1-acc)*100))
60 print('-'*80)
61
62 print(X.shape)
63 print(y.shape)
64
65 y = y.ravel()
66 if y.size > 0 :
67     scores = cross_val_score(bst,X,y, scoring = "accuracy",cv=15)
68     print(" The accuracy by cross validation function is :",scores.mean())
69
```

```
In [ ]: 1 # Seeing the Performance for Xgboost Classifier with default set of paramters
2 # Classifier with DEfault parameters
3
4 bst1 = XGBClassifier()
5
6 # training the Classifier with training set
7 #st1.fit(X_train,y_train ,eval_set=[(X_train, y_train), (X_test, y_test)])
8
9 bst1.fit(X_train,y_train)
10
11 print(" Xgboost Classifier 2 with Default parameters \n",'-'*80)
12
13 print(bst1)
14
15 # predicting the labels for the testing set
16 preds1 = bst1.predict(X_test)
17
18 preds1 = pd.Series(preds1)
19
20 print(y_test.shape)
21 print(preds1.shape)
22
23 preds1 = preds1.values.reshape(-1,1)
24
25 acc_def = round(accuracy_score(y_test, preds1),2)
26
27 print('Accuracy for default classifier by accuracy score function: ', acc_def*100, '%')
28
29 scores = cross_val_score(bst1,X,y,scoring ="accuracy",cv=15)
30 print(" The accuracy by cross validation function is :",round(scores.mean(),2)*100)
31
```



```
In [ ]: 1 # Predicting the Performance by using directly Xgboost Classifier
2 print(y_train.shape)
3
4 # converting the training and test set to sparse matrix by using DMatrix
5 XG_train = xgb.DMatrix(X_train,y_train)
6 XG_test = xgb.DMatrix(X_test,y_test)
7
8
9 params = {
10     'objective':'binary:logistic',
11     'max_depth':4,
12     'silent':1,
13     'eta':1
14 }
15
16 num_rounds = 26
17 watchlist = [(XG_test,'test'), (XG_train,'train')]
18
19 bst2 = xgb.train(params, XG_train, num_rounds,watchlist)
20 #bst = xgb.train(params, XG_train, num_rounds)
21
22 print(bst2)
23
24 preds_prob = bst2.predict(XG_test)
25
26 print(preds_prob)
27
28 print("\nTrain possible labels: ")
29 print(np.unique(XG_train.get_label()))
30
31 print("\nTest possible labels: ")
32 print(np.unique(XG_test.get_label()))
33
34 labels = XG_test.get_label()
35 preds = preds_prob > 0.5 # threshold
36
37 correct = 0
38
39 for i in range(len(preds)):
40     if (labels[i] == preds[i]):
41         correct += 1
```

```
42
43 acc = accuracy_score(labels,preds)
44
45 print('\n Predicted correctly: {0}/{1}'.format(correct, len(preds)))
46 print('Accuracy calculated manually : ', correct*100/len(preds), '%')
47 print('Accuracy by accuracy score function: ', acc*100, '%')
48 print('Error: {0:.4f}'.format(1-correct/len(preds)))
49
50
51 print(" The ROC score is : %.3f "%roc_auc_score(labels,preds))
52 print(" The confusion matrix is :\n")
53 print(confusion_matrix(labels,preds))
54
```

```
In [ ]: 1 # finding the Kfold and Cross Val score for this algorithm
2
3 from sklearn.model_selection import KFold,cross_val_score
4 cv_results = xgb.cv(params, XG_train, num_rounds,nfold=5,metrics="auc")
5 print(cv_results.head())
6 print('-'*80)
7 print(cv_results.tail(1))
8 print(cv_results.mean())
9
```