KNN,SVM & RANDOM FOREST

Final Project Project

Course: CS-634

Name: Anusha Syed

Email: as445@njit.edu.

Libraries used:

numpy, random and pandas for reading files and display of data in data frame.

- Imported KNeighborsClassifier from sklearn.neighbors
- Imported StandardScaler , MinMaxScaler from sklearn.preprocessing
- Imported train test split from sklearn.model selection
- Imported classification_report , confusion_matrix, accuracy_score, f1_score from sklearn.metrics
- Imported LabelEncoder, StandardScaler from sklearn.preprocessing
- Imported RandomForestClassifier from sklearn.ensemble
- Imported StratifiedKFold, RandomizedSearchCV from sklearn.model_selection

→ Used all these libraries for classification according to given model

```
# Importing the important libraries

import numpy as np
import random
import pandas as pd

from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import StandardScaler , MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report , confusion_matrix, accuracy_score, f1_score

from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import StratifiedKFold, RandomizedSearchCV

import warnings
warnings.filterwarnings('ignore')
```

File (dataset) Input:

Here, I have loaded my data set using pandas library and displayed it.

- → df.shape tells how many rows and columns are there
- → df.isnull().sum() tells if there is any null value in data and if it is then given the sum of the number of items that has null values
- → df.describe() shows the actual data stored in your csv file

```
# Loading the dataset

df = pd.read_csv('heart.csv')

df.head()

C* Age Sex ChestPainType RestingBP Cholesterol FastingBS RestingECG MaxHR ExerciseAngina Oldpeak ST Slope HeartDisease
```

>		Age	Sex	ChestPainType	RestingBP	Cholesterol	FastingBS	RestingECG	MaxHR	ExerciseAngina	Oldpeak	ST_Slope	HeartDisea	se
	0	40	М	ATA	140	289	0	Normal	172	N	0.0	Up		0
	1	49	F	NAP	160	180	0	Normal	156	N	1.0	Flat		1
	2	37	М	ATA	130	283	0	ST	98	N	0.0	Up		0
	3	48	F	ASY	138	214	0	Normal	108	Υ	1.5	Flat		1
	4	54	М	NAP	150	195	0	Normal	122	N	0,0	tivate M	lindows	0

```
print(df.shape)
print(df["HeartDisease"].value_counts())
```

C→ (918, 12)
 1 508
 0 410
 Name: HeartDisease, dtype: int64

[53] df.isnull().sum()

Age 0
Sex 0
ChestPainType 0
RestingBP 0
Cholesterol 0
FastingBS 0
RestingECG 0
MaxHR 0
ExerciseAngina 0
Oldpeak 0
ST_Slope 0
HeartD1sease 0
dtype: int64

[54] df.describe()

D		Age	RestingBP	Cholesterol	FastingBS	MaxHR	Oldpeak	HeartDisease
	count	918.000000	918.000000	918.000000	918.000000	918.000000	918.000000	918.000000
	mean	53.510893	132.396514	198.799564	0.233115	136.809368	0.887364	0.553377
	std	9.432617	18.514154	109.384145	0.423046	25.460334	1.066570	0.497414
	min	28.000000	0.000000	0.000000	0.000000	60.000000	-2.600000	0.000000
	25%	47.000000	120.000000	173.250000	0.000000	120.000000	0.000000	0.000000
	50%	54.000000	130.000000	223.000000	0.000000	138.000000	0.600000	1.000000
	75%	60.000000	140.000000	267.000000	0.000000	156.000000	1.500000	1.000000
	max	77.000000	200.000000	603.000000	1.000000	202.000000	6.200000	1.000000

```
[55] df.info()
      <class 'pandas.core.frame.DataFrame'>
      RangeIndex: 918 entries, 0 to 917
      Data columns (total 12 columns):
      # Column Non-Null Count Dtype
                           918 non-null int64
      0 Age
      1 Sex
      1 Sex 918 non-null object
2 ChestPainType 918 non-null object
      3 RestingBP 918 non-null int64
      4 Cholesterol 918 non-null int64
      5 FastingBS 918 non-null int64
6 RestingECG 918 non-null object
7 MaxHR 918 non-null int64
      8 ExerciseAngina 918 non-null object
      9 Oldpeak 918 non-null float64
10 ST_Slope 918 non-null object
11 HeartDisease 918 non-null int64
      dtypes: float64(1), int64(6), object(5)
      memory usage: 86.2+ KB
```

Separating the dependent and independent columns

Here, X is my training data and Y is my label.

```
[56] # Seperating the dependent and the independent columns

X = df.drop("HeartDisease" , axis = 1)
y = df['HeartDisease']
```

Separating the numerical and the categorical columns

Here, cat_cols brings the categorical data like those that don't have any numerical value and num_cols prints those labels that has numerical values to it.

```
[58] print(cat_cols , num_cols)

['Sex', 'ChestPainType', 'RestingECG', 'ExerciseAngina', 'ST_Slope'] ['Age', 'RestingBP', 'Cholesterol', 'FastingBS', 'MaxHR', 'Oldpeak']
```

```
# Seperating the numerical and the categorical columns

cat_cols = list(df.select_dtypes(include=['object']).columns)

num_cols = []

for i in list(X.columns):
    if i not in cat_cols:
        num_cols.append(i)
```

Getting mean and percentile of data

Separating the test and training data set

Here, test_size=0.2 means that 20% is the testing data and rest 80% is the training data.

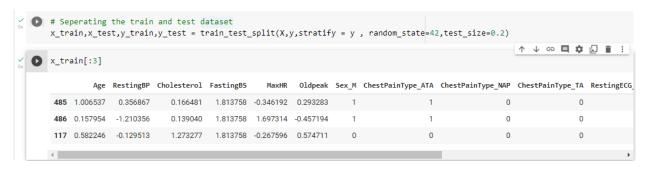
```
[61] # Seperating the train and test dataset
    x_train,x_test,y_train,y_test = train_test_split(X,y,stratify = y , random_state=42,test_size=0.2)
```

Standardadising the features of training and testing dataset

```
[63] # Standardadising the features of training and testing dataset

scaler = StandardScaler()
scaler.fit(X[num_cols])
X[num_cols] = scaler.transform(X[num_cols])
```

Converted categorical data also into numerical data and now again separating the test and training data



KNN Algorithm:

▼ KNN Algorithm

```
[67] knn = KNeighborsClassifier(n_neighbors=7)
    knn.fit(x train,y train)
    y_pred_knn = knn.predict(x_test)
    print(knn.score(x_test,y_test))
    print(classification_report(y_test,y_pred_knn))
    print(confusion_matrix(y_test,y_pred_knn))
    #compare with table
    0.8967391304347826
                precision recall f1-score support
                   0.90
                           0.87 0.88
                                              82
                   0.90
                           0.92
                                    0.91
                                              102
       accuracy
                                     0.90
                                              184
                 0.90 0.89
0.90 0.90
                                   0.90
                                              184
       macro avg
                                    0.90
    weighted avg
                                              184
    [[71 11]
     [ 8 94]]
```

I have used the built in functions extracting from library for this algorithm to set up the training data and labels and printing out the report that shows performance metrics. It also prints the accuracy of the algorithm and prints out the confusion matrix.

Refer to this diagram to interpret results for confusion matrix:

Positive (1) Negative (0) Positive (1) TP FP Negative (0) FN TN

Therefore, 71 is the TP, 11 is FP, 8 is FN and 94 is TN.

SVM Algorithm:

→ SVM Algorithm

```
[68] from sklearn.svm import SVC
      classifier = SVC(kernel='rbf', random_state=27)
      classifier.fit(x_train, y_train)
      SVC(random state=27)
(69) y_pred_svm = classifier.predict(x_test)
 [70] print(classifier.score(x test,y test))
      print(classification_report(y_test, y_pred_svm))
      print(confusion_matrix(y_test,y_pred_svm))
      0.8858695652173914
                 precision recall f1-score support
                    0.89
                            0.85 0.87
                                                82
                      0.89
                             0.91
                                       0.90
                                                102
                                       0.89
                                                184
         accuracy
                    0.89 0.88 0.88
                                                184
         macro avg
      weighted avg
                      0.89
                             0.89
                                     0.89
                                                184
      [[70 12]
       [ 9 93]]
```

I have used the built in functions extracting from library for this algorithm to set up the training data and labels and printing out the report that shows performance metrics. It also prints the accuracy of the algorithm and prints out the confusion matrix.

Refer to this diagram to interpret results for confusion matrix:

Actual Values

		Positive (1)	Negative (0)
d Values	Positive (1)	TP	FP
Predicted	Negative (0)	FN	TN

Therefore, 70 is the TP, 12 is FP, 9 is FN and 93 is TN.

Random Forest Algorithm:

▼ Random Forest Algorithm

```
[71] randomforest = RandomForestClassifier()
      randomforest.fit(x_train,y_train)
      y_pred_rf = randomforest.predict(x_test)
[72] print('Train Accuracy %s' % round(accuracy_score(y_test, y_pred_rf),2))
      print('Train F1-score %s' % f1_score(y_test, y_pred_rf, average=None))
      print(classification_report(y_test, y_pred_rf))
      print("Confusion Matrix:\n",confusion_matrix(y_test, y_pred_rf))
      Train Accuracy 0.86
      Train F1-score [0.84146341 0.87254902]
                  precision recall f1-score support
                      0.84
                              0.84
                                       0.84
                                                   82
                      0.87
                              0.87
                                       0.87
                                                  102
                                                184
184
                                        0.86
          accuracy
      macro avg 0.86 0.86 0.86
weighted avg 0.86 0.86 0.86
                                                 184
      Confusion Matrix:
       [[69 13]
       [13 89]]
```

I have used the built in functions extracting from library for this algorithm to set up the training data and labels and printing out the report that shows performance metrics. It also prints the accuracy of the algorithm and prints out the confusion matrix.

Refer to this diagram to interpret results for confusion matrix:

Positive (1) Negative (0) Positive (1) TP FP Negative (0) FN TN

Therefore, 69 is the TP, 13 is FP, 13 is FN and 89 is TN.

10-folds execution Method:

Declaring list for accuracy, tss, precision:

```
knn_acc = []
knn_tss = []
knn_prec = []
knn_tn = []
knn_tp = []
knn_fn = []
knn_fp = []
svm_acc = []
svm_tss = []
svm_prec = []
svm_tn = []
svm_tp = []
svm_fn = []
svm_fp = []
randmf_acc = []
randmf_tss = []
randmf_prec = []
randmf_tn = []
randmf_tp = []
randmf_fn = []
randmf_fp = []
```

```
for i in range(0,11):
  x_train,x_test,y_train,y_test = train_test_split(X,y,stratify = y , random_state=42,test_size=0.3)
  ## Running KNN Algorithm
  knn.fit(x_train,y_train)
  y_pred_knn = knn.predict(x_test)
  tn, fp, fn, tp = confusion_matrix(y_test, y_pred_knn).ravel()
  knn tn.append(tn)
  knn_tp.append(fp)
  knn_fn.append(fn)
  knn_fp.append(tp)
  acck = (tp + tn) / (tn + fp + fn + tp)
  tss = (tp / (tp - fn)) - (fp / (fp + tn))
  precision = tp / (tp + fp)
  knn_acc.append(acck)
  knn_tss.append(tss)
  knn_prec.append(precision)
```

```
## Running Random Forest Algorithm
randomforest.fit(x_train,y_train)
y_pred_rf = randomforest.predict(x_test)

tn, fp, fn, tp = confusion_matrix(y_test, y_pred_rf).ravel()
randmf_tn.append(tn)
randmf_tp.append(fp)
randmf_fn.append(fn)
randmf_fp.append(tp)

accrf = (tp + tn) / (tn + fp + fn + tp)
tss = (tp / (tp - fn)) - (fp / (fp + tn))
precision = tp / (tp + fp)

randmf_acc.append(accrf)
randmf_tss.append(tss)
randmf_prec.append(precision)
```

```
## Running SVM Algorithm
classifier.fit(x_train, y_train)
y_pred_svm = classifier.predict(x_test)

tn, fp, fn, tp = confusion_matrix(y_test, y_pred_svm).ravel()
svm_tn.append(tn)
svm_tp.append(fp)
svm_fn.append(fn)
svm_fp.append(tp)

accs = (tp + tn) / (tn + fp + fn + tp)
tss = (tp / (tp - fn)) - (fp / (fp + tn))
precision = tp / (tp + fp)

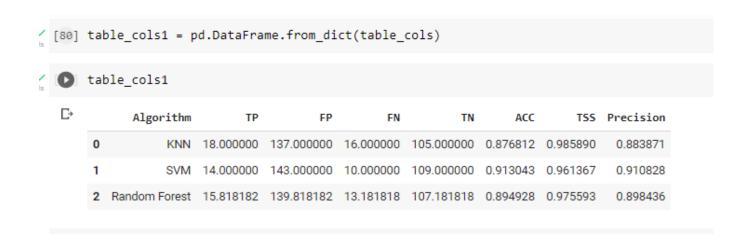
svm_acc.append(accs)
svm_tss.append(tss)
svm_prec.append(precision)
```

Average:

```
## Average:
        avg knn acc = sum(knn acc) / len(knn acc)
        avg_knn_tss = sum(knn_tss) / len(knn_tss)
        avg knn prec = sum(knn prec) / len(knn prec)
        avg knn tn = sum(knn tn) / len(knn tn)
        avg knn tp = sum(knn tp) / len(knn tp)
        avg knn fn = sum(knn fn) / len(knn fn)
        avg knn fp = sum(knn fp) / len(knn fp)
/ [75] ## Average:
       avg svm acc = sum(svm acc) / len(svm acc)
       avg svm tss = sum(svm tss) / len(svm tss)
       avg svm prec = sum(svm prec) / len(svm prec)
       avg svm tn = sum(svm tn) / len(svm tn)
       avg svm tp = sum(svm tp) / len(svm tp)
       avg_svm_fn = sum(svm_fn) / len(svm_fn)
       avg svm fp = sum(svm fp) / len(svm fp)
    ## Average:
         avg_randmf_acc = sum(randmf_acc) / len(randmf_acc)
         avg randmf tss = sum(randmf tss) / len(randmf tss)
         avg randmf prec = sum(randmf prec) / len(randmf prec)
         avg randmf tn = sum(randmf tn) / len(randmf tn)
         avg randmf tp = sum(randmf tp) / len(randmf tp)
         avg randmf fn = sum(randmf fn) / len(randmf fn)
         avg randmf fp = sum(randmf fp) / len(randmf fp)
y [o] table_cols = {'Algorithm' : ['KNN' , 'SVM' , 'Random Forest'] , 'TP' : [avg_knn_tp , avg_svm_tp , avg_randmf_tp] , 'FP' : [ε
[77] 'FN': [avg_knn_fn, avg_randmf_fn], 'TN': [avg_knn_tn, avg_svm_tn, avg_randmf_tn], 'ACC': [avg_knn_acc, avg
🐧 :c], 'TSS' : [avg_knn_tss , avg_svm_tss, avg_randmf_tss] , 'Precision' : [avg_knn_prec, avg_svm_prec, avg_randmf_prec]}
```

Table:

This is the final table that shows accuracy of algorithm after running it ten times.



Result:

Which algorithm performs better and why? Justify your ansIr

- → Accuracy for KNN was coming 90 and after 10 folds it is coming as 87
- → Accuracy for SVM was coming 89 and after 10 folds it is coming as 91
- → Accuracy for Random Forest was coming 86 and after 10 folds it is coming as 89

SVM algorithm performs better and you can see that after training the algorithm 10 times the accuracy of it is better. It is also memory efficient. It also performs faster prediction.